# LEVEL II

# NAVAL POSTGRADUATE SCHOOL
## Monterey, California

DTIC
SELECTE
JUL 1 5 1980
S          D
C

# THESIS

CLASSIFICATION TECHNIQUES FOR
MULTIVARIATE DATA ANALYSIS

by

Jin Ki Lee

March 1980

Thesis Advisor:          F. R. Richards

Approved for public release; distribution unlimited.

80 7 14 108

| REPORT DOCUMENTATION PAGE | | READ INSTRUCTIONS BEFORE COMPLETING FORM |
|---|---|---|
| 1. REPORT NUMBER | 2. GOVT ACCESSION NO. AD-A086521 | 3. RECIPIENT'S CATALOG NUMBER |
| 4. TITLE (and Subtitle) Classification Techniques for Multivariate Data Analysis. | | 5. TYPE OF REPORT & PERIOD COVERED Master's Thesis; March 1980 |
| | | 6. PERFORMING ORG. REPORT NUMBER |
| 7. AUTHOR(s) Jin Ki Lee | | 8. CONTRACT OR GRANT NUMBER(s) |
| 9. PERFORMING ORGANIZATION NAME AND ADDRESS Naval Postgraduate School Monterey, California 93940 | | 10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS |
| 11. CONTROLLING OFFICE NAME AND ADDRESS Naval Postgraduate School Monterey, California 93940 | | 12. REPORT DATE 28 March 1980 |
| | | 13. NUMBER OF PAGES 120 |
| 14. MONITORING AGENCY NAME & ADDRESS(If different from Controlling Office) Naval Postgraduate School Monterey, California 93940 | | 15. SECURITY CLASS. (of this report) Unclassified |
| | | 15a. DECLASSIFICATION/DOWNGRADING SCHEDULE |

16. DISTRIBUTION STATEMENT (of this Report)

Approved for public release; distribution unlimited.

17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, If different from Report)

18. SUPPLEMENTARY NOTES

251450

19. KEY WORDS (Continue on reverse side if necessary and identify by block number)

multivariate analysis, principal components analysis, variance-covariance matrix, correlation, lagranian technique discriminant analysis, distance measure (metric), hierarchical clustering, nonhierarchical clustering, similarity matrix

20. ABSTRACT (Continue on reverse side if necessary and identify by block number)

The multivariate analysis techniques of cluster analysis, principal components analysis, and discriminant analysis are examined in this thesis. The theory and applications of each of the techniques are discussed. Computer software available at the Naval Postgraduate School is discussed and sample jobs are included.

A hierarchical cluster analysis algorithm, available in the IMSL software package, is applied to a set of data extracted from a

group of subjects for the purpose of partitioning a collection of 26 attributes of a weapon system into six clusters of super-attributes.

A nonhierarchical clustering procedure, principal components analysis, and discriminant analysis were all applied to a collection of data on tanks considering of twenty-four observations of ten attributes of tanks. The cluster analysis shows that the tanks cluster somewhat naturally by nationality. The principal components analysis and the discriminant analysis show that tank weight is the single most important discriminator among nationality.

| Accession For | | |
|---|---|---|
| NTIS GRA&I | | ☑ |
| DDC TAB | | ☐ |
| Unannounced | | ☐ |
| Justification | | |
| By | | |
| Distribution/ | | |
| Availability Codes | | |
| Dist | Avail and/or Special | |
| A | | |

Classification Techniques for
Multivariate Data Analysis

by

Jin Ki Lee
Lieutenant Colonel, Korean Army
B.S., Korean Military Academy, 1966

Submitted in partial fulfillment of the
requirements for the degree of

MASTER OF SCIENCE IN OPERATIONS RESEARCH

from the

NAVAL POSTGRADUATE SCHOOL
March 1980

Author _____

Approved by: _____
Thesis Advisor

_____
Co-Advisor or Second Reader

_____
Chairman, Department of Operations Research

_____
Dean of Information and Policy Sciences

3

## ABSTRACT

The multivariate analysis techniques of cluster analysis, principal components analysis, and discriminant analysis are examined in this thesis. The theory and applications of each of the techniques are discussed. Computer software available at the Naval Postgraduate School is discussed and sample jobs are included.

A hierarchical cluster analysis algorithm, available in the IMSL software package, is applied to a set of data extracted from a group of subjects for the purpose of partitioning a collection of 26 attributes of a weapon system into six clusters of superattributes.

A nonhierarchical clustering procedure, principal components analysis, and discriminant analysis were all applied to a collection of data on tanks considering of twenty-four observations of ten attributes of tanks. The cluster analysis shows that the tanks cluster somewhat naturally by nationality. The principal components analysis and the discriminant analysis show that tank weight is the single most important discriminator among nationality.

TABLE OF CONTENTS

5

## LIST OF TABLES

# LIST OF FIGURES

# I.   DISCUSSION OF MULTIVARIATE DATA ANALYSIS

## A.   INTRODUCTION

As a set of statistical techniques, multivariate data analysis is concerned with data collected on several dimensions of the same observations.  Techniques can be used for many purpose in the behavioral, mathematical, and administrative sciences - ranging from rigidly controlled experiments to explain relationships assumed to be present in a large mass of data to attempts to cluster similar elements or to find functions of the variables that will best discriminate among preselected subpopulations of the observations.

The heart of any multivariate analysis consists of the data matrix.  This matrix is a table that gives the results of a number of observations on a number of variables simultaneously (Table I).

Illustrative Data Matrix

| Observations | Variables | | | | |
|---|---|---|---|---|---|
| | 1 | 2 | 3 .... | j .... | p |
| 1 | $x_{11}$ | $x_{12}$ | $x_{13}$ | $x_{1j}$ | $x_{1p}$ |
| 2 | $x_{21}$ | $x_{22}$ | $x_{23}$ | $x_{2j}$ | $x_{2p}$ |
| . | . | . | . | . | . |
| . | . | . | . | . | . |
| . | . | . | . | . | . |
| i | $x_{i1}$ | $x_{i2}$ | $x_{i3}$ | $x_{1j}$ | $x_{ip}$ |
| . | . | . | . | . | . |
| . | . | . | . | . | . |
| . | . | . | . | . | . |
| n | $x_{n1}$ | $x_{n2}$ | $x_{n3}$ | $x_{nj}$ | $x_{np}$ |

TABLE I.

The table consists of a set of observations (the n rows) and a set of measurements on those observations (the p columns). Cell entries represent the value $x_{ij}$ of observation i on variable j . The values are characteristics of the observations and serve to define the observations in any specific study. The cell values may consist of nominal, ordinal, interval, or ratio-scaled measurements, or varicus combinations of these across columns.

In a general sense "multivariate" analysis would concern two main features:

1. The multivariate character lies in the multiplicity of the $p$ variables, not in the size of the set $n$.

2. The variables are dependent among themselves so that we can not split off one or more from the others and consider it by itself. The variables must be considered together.

There are three characteristics often used as a basis for the classification of multivariate analysis:

1. whether one's principal focus is on the objects or on the variables of the data matrix;

2. whether the data matrix is partitioned into criterion and independent subsets, and the number of variables in each;

3. whether the cell values represent nominal, ordinal, or interval scale measurements.

This classification results in four major subdivisions of interest:

1. single criterion, multiple predictor association, including multiple regression, analysis of variance and covariance, and two-group discriminant analysis;

2. multiple criterion, multiple predictor
   association, including canonical correla-
   tion, multivariate analysis of variance
   and covariance, and multiple discriminant
   analysis;

3. analysis of variable interdependence,
   including factor analysis, multidimen-
   sional scaling, and other types of
   dimension-reducing methods;

4. analysis of interobject similarity,
   including cluster analysis and other
   types of grouping procedures.

The first two categories involve dependence structures
where the data matrix is partitioned into criterion and
independent subsets; in both cases interest is focused on
the variables. The last two categories are concerned with
interdependence - either focusing on variables or on
observations. Within each of four categories, various
techniques are differentiated in terms of the type of scale
assumed.

In this research, we consider only the following
techniques of multivariate analysis:

1. Principal components analysis

2. Discriminant analysis

3. Cluster analysis

## II. PRINCIPAL COMPONENTS ANALYSIS

The basic idea of principal components analysis is to describe the dispersion of an array of  n  points in p-dimensional space by introducing a new set of orthogonal linear coordinates so that the sample variances of the given data points with respect to these derived coordinates are in decreasing order of magnitude.  Thus the first principal component is such that the projection of given points onto it have maximum variance among all possible linear coordinates; the second principal component has maximum variance subject to being orthogonal to the first; and so on.

Suppose that the random variables  $X_1$ , $X_2$ , ......, $X_p$  of interest have a certain multivariate distribution with finite mean vector  u  and variance-covariance matrix  $\Sigma$ .

From this population a sample of  n  independent observation vectors has been drawn.  The observation can be written as the usual  $n \times p$  data matrix.

$$X = \begin{bmatrix} x_{11} \cdots \cdots \cdots \cdots x_{ip} \\ \vdots \qquad\qquad \vdots \\ \vdots \qquad\qquad \vdots \\ x_{n1} \cdots \cdots \cdots \cdots x_{np} \end{bmatrix} = \begin{bmatrix} X_1' \\ \vdots \\ \vdots \\ X_n' \end{bmatrix} \qquad (1)$$

The estimate of $\Sigma$ will be the usual sample variance-covariance matrix $S$ defined as follows:

$$S = \frac{1}{n-1} A$$

$$A = \sum_{j=1}^{n} (X_j - \bar{X})(X_j - \bar{X})' \qquad (2)$$

The information we shall need for our principal components analysis will be contained in $S$. However, it will be necessary to make a choice of measures of dependence: should we work with the variances and covariances of the observations, and carry out the analysis in original unit of the responses, or would a more accurate picture of the dependence pattern be obtained if each $x_{ij}$ were transformed to a standarized score

$$Z_{ij} = \frac{x_{ij} - \bar{x}_j}{s_j}$$

and the correlation matrix $R$ employed? The components obtained from $S$ and $R$ in general not the same, nor is it possible to pass from one solution to the other by a simple scaling of the coefficients.

If the responses are in widely different units (i.e., number of crew, weight in tons, speed in kilometer per hour, etc.) with large differences in the magnitudes, linear compounds of original quantities would have little

meaning and standarized variates and correlation matrix
should be employed.  Conversely, if the responses are
reasonably commensurable, the covariance form has a greater
statistical appeal, for the i-th principal component is
that linear compound of the responses which explains the
i-th largest portion of the total response variance, and
maximization of such total variance of standard scores is
rather artificial.

The first principal component of the complex of sample
values of the responses $X_1$ , $X_2$ , .........., $X_p$ is the
linear compound

$$Y_1 = a_{11}X_1 + .........+ a_{p1}X_p \qquad (3)$$

whose coefficients $a_{i1}$ are the elements of the eigenvector
associated with the greatest eigenvalue $\lambda_1$ of the sample
variance-covariance matrix of the responses. The $a_{i1}$ are
are unique up to multiplication by a scale factor, and if
they are scaled so that $a'_1 a_1 = 1$ , the eigenvalue $\lambda_1$ is
interpretable as the sample variance of $Y_1$ .

Numerical representation of the first principal compo-
nent is to find the vector $A_1$ such that

$$Y_1 = a_{11}X_1 + .........+ a_{p1}X_p$$
$$= A'_1 X \qquad (4)$$

which maximizes sample variance

$$s_{Y_1}^2 = \sum_{i=1}^{p} \sum_{j=1}^{p} a_{i1} a_{j1} s_{ij}$$

$$= A'_1 S A_1 \tag{5}$$

for all coefficient vectors normalized so that $A'_1 A_1 = 1$. To determine the coefficients, the normalization constraint is introduced by means of Lagrange multiplier and the resulting expression is differentiated with respect to $A'_1$:

$$\frac{\partial}{\partial A_1}[S_{Y_1}^2 - \lambda_1(1 - A_1'A_1)] = \frac{\partial}{\partial A_1}[A_1'SA_1 + \lambda_1(1 - A_1'A_1)]$$

$$= 2(S - \lambda_1 I)A_1 \tag{6}$$

The coefficients must satisfy the $p$ simultaneous linear equations.

$$(S - \lambda_1 I)A_1 = 0 \tag{7}$$

If the solution to these equation is to be other than the null vector, the value of $\lambda_1$ must be chosen so that

$$|S - \lambda_1 I| = 0 \tag{8}$$

$\lambda_1$ is thus an eigenvalue of the variance-covariance matrix, and $A_1$ is its associated eigenvector. To determine which of the $p$ eigenvalues should be used, premultiply the

15

the system of equation (7) by $A_1'$ . Since $A_1'A_1 = 1$ , it follows that

$$\lambda_1 = A_1' \ S \ A_1 = s_{Y_1}^2$$

But the coefficient vecotr was chosen to maximize this variance, and therefore, $\lambda_1$ must be the greatest eigenvalue of S .

The second principal component is that linear compound

$$Y_2 = a_{12}X_1 + \ldots\ldots\ldots + a_{p2} X_p \qquad (9)$$

whose coefficients have been chosen, subject to the constraints

$$
\begin{aligned}
A_2' \ A_2 &= 1 \\
A_1' \ A_2 &= 0
\end{aligned}
\qquad (10)
$$

so that the variance of $Y_2$ , $A_2' \ S \ A_2$ , is a maximum. The first constraint is merely a scaling to assure the uniqueness of the coefficients, while the second requires that $A_1$ and $A_2$ be orthogonal.

The coefficients of the second component can also be found by the Lagrangian technique with two multipliers $\lambda_2$ and $\mu$ . Differentiating this with respect to $A_2$ gives:

16

$$\frac{\partial}{\partial A_2}[A_2' \ S \ A_2 \ + \ \lambda_2(1 \ - \ A_2' \ A_2) \ + \ \mu A_1' \ A_2]$$

$$= \ 2(S \ - \ \lambda_2 I)A_2 \ + \ \mu A_1 \qquad (11)$$

If the right-hand side is set equal to 0 and premultiplied by $A_1'$ , it follows from the normalization and orthogonality conditions that

$$2 \ A_2' \ S \ A_2 \ + \ \mu \ = \ 0 \qquad (12)$$

Similar premultiplication of the equation (7) by $A_2'$ implies that

$$A_2' \ S \ A_2 \ = \ 0 \qquad (13)$$

and hence $\mu \ = \ 0$ . The second vector must satisfy

$$(S \ - \ \lambda_2 I)A_2 \ = \ 0 \qquad (14)$$

And it follows that the coefficients of the second component are thus the elements of the eigenvector corresponding to the second greatest eigenvalue. The remaining principal components are found in their turn in the same manner from the other eigenvectors.

Thus the $j$-th principal component of the sample of p-variate observations is the linear compound

$$Y_j \ = \ a_{1j}X_1 \ + \ \dots\dots\dots + \ a_{pj}X_p \qquad (15)$$

17

whose coefficients are the elements of the eigenvector of
the sample variance-covariance matrix S corresponding to
the j-th largest eigenvalue $\lambda_j$ . If $\lambda_i \neq \lambda_j$ , the
coefficients of the i-th and j-th components are
necessarily orthogonal; if $\lambda_i = \lambda_j$ , the elements can be
chosen to be orthogonal, although an infinity of such
orthogonal vectors exists. The sample variance of the
j-th components is $\lambda_j$ , and the total system variance is
thus

$$\lambda_1 + \lambda_2 + \ldots\ldots\ldots + \lambda_p = tr\ S \qquad (16)$$

The importance of the j-th component in a more parsimonious
description of the system is measured by

$$\frac{\lambda_j}{tr\ S} \qquad (17)$$

which gives the fraction of the total variance contributed
to the j-th component.

# III.  DISCRIMINANT ANALYSIS

## A.  INTRODUCTION

The basic idea of discriminant analysis consists of assigning an individual from a group of individuals to one of several known or unknown distinct propulations, on the basis of observations on several characters of the individual or group and a sample of observations on these characters from the populations if these are unknown.

Fisher (1936) was the first to suggest a linear function of variables representing different characters, hereafter called the linear discriminant function (discriminator) for classifying an individual into one of two populations.  Later research extended the analysis to classification into one of  k  populations.

For the univariate case Fisher suggested a rule which classifies an observation  x  into the  i-th  univariate population if

$$X - \bar{X}_i = \min (X - \bar{X}_1 , X - \bar{X}_2) , \quad i = 1,2 \quad (18)$$

where  $\bar{X}_i$  is the sample mean based on a sample of size  $N_1$  from  i-th  population. For two  p-variate  populations  $\pi_1$  and  $\pi_2$  (with the same covariance matrix) Fisher replaced the vector random variable by an optimum linear combination of its components obtained by maximizing the

ratio of the difference of the expected values of a linear combination under $\pi_1$ and $\pi_2$ to its standard deviation. He then used his univariate discrimination method with this optimum linear combination of components as the random variable.

Rao (1948) considered the problem of classifying people into one of these populations castes of India. He assumed that each of the three populations could be characterized by four variables - structure $(x_1)$, sitting height $(x_2)$, nasal depth $(x_3)$, and nasal height $(x_4)$ - of each member of the population. On the basis of sample observations on these characters from the three populations the problem is to classify an individual with observation $X = (x_1, x_2, x_3, x_4)^T$ into one of three populations. He used a linear discriminator to obtain the solution.

B.  THEORY

In general, the underlying assumptions of discriminant analysis are:

1.  the groups being investigated are discrete and identifiable;

2.  each observation in each group can be described by a set of measurements on $p$ characteristics or variables;

3.  these $p$ variables are assumed to have a multivariate normal distribution in each population.

The purposes of discriminant analysis are:

1. to test for mean group differences and to describe the overlaps among groups;

2. to construct classification schemes based upon the set of p variables in order to assign previously unclassified observations to the appropriate groups.

Hence, the problem of studying the direction of group differences is, equivalently, a problem of finding a linear combination of the original independent variables that shows large differences in group means. In short, discriminant analysis is a method for determining scuh linear combinations.

The first step toward determining a linear combination of a set of variables such that several group means on this linear combination will differ widely among themselves, is to decide on a criterion for measuring such group-mean differences. Once a linear combination has been constructed, that means there is just a single tranformed variable. Hence, the F-ratio for testing the significance of the over all difference among several group means on a single variable suggests an appropriate criterion.

$$F = \frac{v'Bv}{v'Wv} = \lambda \qquad (19)$$

where $v' = (v_1, v_2, \ldots\ldots\ldots\ldots, v_p)$, a set of weight which maximizes $\lambda$ .

$$B = \sum_{i=1}^{G} N_i (\bar{x}_{i.} - \bar{x}_{..})(\bar{x}_{i.} - \bar{x}_{..})'$$

$$W = \sum_{i=1}^{G} \sum_{I=1}^{n_i} (x_{ij} - \bar{x}_{i.})(x_{ij} - \bar{x}_{i.})'$$

$x_{ij}$ is the jth observation vector in the i-th group.

$\bar{x}_{..}$ is the grand mean vector of the data.

$G$ is the number of groups.

$n_i$ is the number of observations in the ith group.

Prime notation indicates transpose.

This ratio $\lambda$, called the discriminant criterion, was originally proposed by Fisher in connection with his two-group discriminant function. Once a criterion for group differentiation has been determined, a set of weights, $(v_1 \; v_2 , \ldots\ldots , v_p)$, which maximizes this criterion, should be determined. This is accomplished by taking the partial derivative of $\lambda$ with respect to each component $v_i$ of $v$ and setting the result equal to zero.

$$\frac{\partial \lambda}{\partial v} = \frac{2[(Bv)(v'Wv) - (v'Wv)(Wv)]}{(v'Wv)^2}$$

$$= \frac{2(Bv - Wv)}{v'Wv} = 0$$

(20)

which is equivalent to

$$(B - \lambda W)v = 0$$
$$(W^{-1}B - \lambda I)v = 0 \qquad (21)$$

This equation is of the form

$$(A - \lambda I)v = 0 \qquad (22)$$

It's solution, yielding the eigenvalues $\lambda_p$ and associated
eigenvectors $V_p$ of the matrix $A$, is therefore the same
as in the principal components analysis, and thus the solved
problem satisfies the problem of maximizing the discriminant
criterion.

In the last equation, the number of non-zero eigenvalues
of a square matrix $A$ is equal to the rank of $A$. With
$W^{-1}B$ playing the role of $A$, the number of non-zero eigen-
values depends on the rank of $B$, since the rank of the
product of two matrices can not exceed the smaller of the two
factor matrices' ranks, and $W^{-1}$ (being nonsingular) must be
of full rank $p$, while the rank of $B$ is usually smaller
than $p$. Thus it is possible to denote the rank of $B$ by
$r = \min (G-1,p)$.

From the fact that the eigenvalues $\lambda_p$ are the values
assumed by the discriminant criterion for linear combination
using the elements of the corresponding eigenvectors $P$
as combining weights, it is clear that the eigenvector

$V_1' = (v_{11}, v_{12}, \ldots \ldots \ldots, v_{1p})$ provides a set of weights such that the transformed variable

$$Y_1 = v_{11}X_1 + v_{12}X_2 + \ldots \ldots \ldots + v_{1p}X_p \qquad (23)$$

has the largest discriminant-criterion, $\lambda$ , achievable by any linear combination of the $p$ independent variables.

What are the properties of the remaining eigenvectors, $v_2, v_3, \ldots \ldots, v_p$ ? The second discriminant function $Y_2 = v_{21}X_1 + v_{22}X_2 + \ldots \ldots + v_{2p}X_p$ whose weights are the elements of the eigenvector $v_2$ associated with the second largest eigenvalue $\lambda_2$ of $W^{-1}B$ has the largest discriminant-criterion among those linear combinations of the $X_i$ that are uncorrelated with the first discriminant function in the total sample observation. Its proof is analogous of that of princpal components analysis. Each discriminant function has a relative (or conditional) maximum value for its discriminant criterion. Therefore, it needs nonly to show that $Y_2$ is uncorrelated with $Y_1$ . Noting that this correlation is proportional to $v_1'Tv_2$ (where $T = W + B$), we have to prove that $v_1'Tv_2 = 0$ .

$$(B - \lambda_i W)v_i = 0 \qquad \text{for each i} \qquad (24)$$

hence,

$$Bv_1 = \lambda_i\, Wv_1 \quad \text{and} \quad Bv_2 = \lambda_2 Wv_2$$

24

premultiplying these equations by $v_2'$ and $v_1'$ respectively,

$$v_2'Bv_1 = \lambda_1 v_2'Wv_1$$

$$v_1'Bv_2 = \lambda_2 v_1'Wv_2 \tag{25}$$

taking the transpose of both sides of the first equation (B and W are symmetric)

$$v_1'BV_2 = \lambda_1 v_1'WV_2$$

thus

$$\lambda_1 V_1'WV_2 = \lambda_2 V_1'WV_2$$

$$(\lambda_1 - \lambda_s)V_1'WV_2 = 0$$

since $\quad\quad\quad \lambda_1 \neq \lambda_2, \; V_1'WV_2 = 0$

therefore, $V_1'WV_2 = 0$ which means $V_1$ and $V_2$ are uncorrelated, and $Y_2$ has this property: its discriminant-criterion value, $\lambda_2$ , is the largest achievable by any linear combination of X's that is uncorrelated (in the total sample) with $Y_1$ . Similarly

$$Y_3 = v_{31}X_1 + v_{32}X_2 + \ldots\ldots + v_{3p}X_p \tag{36}$$

has the largest possible discriminant-criterion value $(\lambda_3)$ among all linear combinations of the X's that are uncorrelated with both $Y_1$ and $Y_2$; and so on until $Y_r$ using the

elements of $V_r$ as weights, has the largest possible discriminant-criterion value among linear combinations that are uncorrelated with all the preceding linear combinations $Y_1, Y_2, \ldots\ldots, Y_{r-1}$ . The linear combinations $Y_1, Y_2, \ldots\ldots, Y_r$ are called the first, second, $\ldots\ldots$, r th (linear) discriminant functions for optimally differentiating among the g given groups.

The situation here is reminiscent of principal components analysis. There, the dimension corresponding to the first component had maximum variance; the second-component dimension had maximum variance among those uncorrelated with the first; and so on. In discriminant analysis, the ratio of between-to within-groups sums-of-squares merely takes the place of variance as the criterion in determining the successive dimensions. However, an important difference between the dimensions identified in discriminant analysis and those in component analysis is that the former are generally not mutually orthogonal in the test space, even though they are uncorrelated. That is, the axis representing the discriminant functions are not a subset of axes obtainable by rigid rotation of the original system of p axes; the discriminant rotation in an oblique rotation.

Just as in the principal components analysis, the dimensions represented by the discriminant functions may be interpreted meaningfully. Even if they are not, it may be possible to achieve parsimony by reducing the dimensionality of the space needed to describe group differences. In

26

seeking to interpret the discriminant functions, the goal is to determine which of the original  p  variables contribute most to each function.  For this prupose, comparison of the realtive magnitudes of the combining weights as given by the elements of each eigenvector of  $W^{-1}B$  is inappropriate because these are weights to be applied to the variables in raw-score scales, and are hence affected by the particular unit used for each variable.

To eleminate the spurious effects of units of measurement on the magnitudes of combining weights, standarized variables should be used.

The relative magnitudes of these standarized weights may be assessed by multiplying each raw-score weight by the standard deviation of the corresponding variable as computed from the within-groups SSCP (Sum of Squares, Cross product) matrix.  This amounts to multiplying each element of a given eigenvector  $V_m$  by the square root of the corresponding diagonal element of  W .  Thus, for each  m , define

$$v_{mi}{}^* - w_{ii} v_{mi} \qquad i = 1,2,\ldots\ldots,p \qquad (27)$$

as the standarized discriminant weights.  The relative contribution of the  i th  variable to the  m th  discriminant function may then be gauged by the magnitude of  $v_{mi}{}^*$  in comparison with the other weights  $v_{mj}{}^*$ .

Up to this point, it has been shown that the dimensionality of the discriminant space is equal to the number of

nonzero eigenvalues of $W^{-1}B$ , which is the smaller of the two numbers, $G-1$ and $p$ . It may often happen, that the number of significant discriminant dimensions may be even smaller. That is, not all of the discriminant function may represent dimensions along which statistically significant group differences occur.

## C.  SIGNIFICANCE TEST IN DISCRIMINANT ANALYSIS

A basic quantity in testing the significance of the overall difference among several group centroids (mean vectors) the ratio of the determinants of the within-groups and the total SSCP matrices, known as Wilks' $\Lambda$ criterion.

$$\Lambda = \frac{|W|}{|T|} \tag{28}$$

Motivation for use of this equation may be seen as follows:

$$\frac{1}{\Lambda} = \frac{|T|}{|W|} \quad |W^{-1}T|$$

$$= |W^{-1}(W + B)| \tag{29}$$

$$= (1 + \lambda)(1 + \lambda_2); \ldots, (1 + \lambda_r)$$

where $\lambda_1, \lambda_2, \ldots, \lambda_r$ are the nonzero eigenvalues of $W^{-1}B$ . Consequently, Bartlet's $V$ statistic for testing the significance of an observed value can be expressed as

28

$$V = - [N - 1 - (p + G)/2]\ln\Lambda$$

$$= [N - 1 - (p + G)/2] \ln [(1 + \lambda_1)(1 + \lambda_r)] \qquad (30)$$

$$= [N - 1 - (p + G)/2] \sum_{m=1}^{r} \ln(1 + \lambda m)$$

This statistic is distributed approximately chi-square with $p(G-1)$ degrees of freedom.

Because of the uncorrelatedness of the successive discriminant functions, the successive terms $\ln(1 + \lambda_m)$ in the last expression above are statistically independent (assuming multivariate normality of the original $p$ variables). As a result, the additive components of $V$ are each approximately distributed as a chi-square variate. More specifically, the $m$ th component,

$$V_m = [N - 1 - (p + G)/2] \ln (1 + \lambda_m) \qquad (31)$$

is approximately chi-square with $p \pm G - 2m$ degrees of freedom. It may be readily verified that the sum of the number of degree of freedom (n.d.f) of the $r$ components, that is, $(p + G - 2) + (p + G = 4) + \ldots\ldots+(p + G - 2r)$ , is equal to $p(G - 1)$ regardless of whether $r = G - 1$ or $p$ .

Consequently, when we cumulatively subtract $V_1, V_2$ , and so on from $V$ , the remainder each time is also a chi-square variate; and these successive remainders become appropriate statistics for testing whether the residual discrimination after removing the first discriminant

29

function, the first and second discriminant function, and so forth, is statistically significant. The successive test statistics and their n.d.f.'s may be summarized as follows:

| Residual After Removing | Approximate - Statistic | n.d.f. |
|---|---|---|
| First discriminant Function | $V - V_1$ | $p(G-1) - (p+G-2)$ $= (p-1)(G-2)$ |
| First 2 discriminant Function | $V - V_1 - V_2$ | $(p-1)(G-2)-(p+G-4)$ $=(p-2)(G-3)$ |
| First 3 discriminant Function | $V - V_1 - V_2 - V_3$ | $(p-2)(G-3)-(p+G-6)$ $=(p-3)(G-4)$ |
| ⋮ | ⋮ | ⋮ |
| First s discriminant Function | $V-V_1-V_2-V_3\ldots-V_s$ | $(p-s)(G-(s+1))$ |

As soon as the residual, after removing the first s discriminant functions becomes smaller than the prescribed percentile point (that is, the $100(1 - \alpha)$th percentile) of the appropriate chi-square distribution, we may conclude that only the first s discriminant functions are significant at that $\alpha$ level. If the number of significant discriminant functions thus found is smaller than r (as will often be the case), we will have effected a further reduction in the dimensionality of the space required to describe the differences among the G groups from which

our sample groups were drawn. The remaining r-s dimensions may be regarded as immaterial for population differentiation, since our sample differences along these dimensions can be attributed to sampling error.

# IV.  CLUSTER ANALYSIS

## A.  ORIGIN AND THEORY

Clustering is the grouping of similar objects.  The
principal functions of clustering are to name, to display,
to summarize, to predict, and to aid in interpretation of
data with many dimensions.  Clustering techniques were
first developed in the field of biological taxonomy.  It is
one of several methodologies included in the broader cate-
gory called classification.

The cluster analysis problem is the last step we
consider in the progression of category sorting problems.
While in discriminant analysis some part of the structure
is known and missing information is estimated from labeled
samples, the operational objectives of clustering is to
classify new observations, that is, recognize them as members
of one category or another.  In cluster analysis little or
nothing is known about the category structure.  All that is
available is a collection of observations whose category
membership are known.  We seek to discover a category
structure which fits the observations.  The problem may be
stated as one of finding the "natural groups", which means
to sort the observations into groups such that the degree of
"natural association" is high among members of the same
group and low between members of different groups.

Cluster analysis techniques have been applied in many fields of study. The literature is both voluminous and diverse, the terminology differing from one field to another. "Numerical taxonomy" is frequently substituted for cluster analysis among biologists, botanists, and ecologists, while some social scientists may refer "typology". Other frequently encountered terms are pattern recognition and partitioning. While discriminant analysis has been studied by statisticians for nearly 45 years, cluster analysis has only recently come to statistical notice. Any method which partition a set of objects into subsets on the basis of measurements taken on every object qualifies as a clustering method.

Most of the well known clustering techniques fall into one of two main categories: (1) hierachical and (2) non-hierachical (partitioning). The former is one in which every cluster obtained at any stage is a merger of clusters at previous stages. The nonhierachial procedures however form new clusters by lumping and splitting old ones. We consider both categories shortly.

In a geometric sense, every observation may be viewed as a point in p-dimensional Euclidean space. This swarm of data points may contain dense regions or "clouds" of data points which are separable from other regions containing a low density of points. These denser regions constitute what are known as clusters. In one and two dimensional cases, it is easy to visualize and to detect the clusters from scatter

33

plots, assuming that the clusters exist. In higher dimensions, clustering becomes extremely difficult without the aid of a computer.

Mathematical clustering techniques usually require a measure of similarity to be defined for every pairwise combination of the entities to be clustered. In order to solve the cluster problem, it is desirable to define the terms "similarity" and "difference" in a quantitative fashion. A researcher would assign two observations to the same group if the distance between them is sufficiently small, or to different clusters if this distance is sufficiently large.

At this point, two questions may be brought on. The first one is "how do we measure the distance between the observations?" and the second one is "how small is small enough?" and how large is large enough? These will be discussed in the following sections.

## B. MEASURES OF DISTANCE

### 1. General

Let $E_p$ be a symbolic representation for a measurement in p-dimensional space and let X,Y, and Z be any of these points in $E_p$. Then any nonnegative real-valued function $D(X,Y)$ satisfying the following conditions qualifies as a distance function (or metric).

1. $D(X,Y) = 0$      if and only if $X = Y$

2. $D(X,Y) \geq 0$      for all X and Y in $E_p$

34

3. $D(X,Y) = D(Y,X)$

4. $D(X,Y) \leq D(X,Z) + D(Y,Z)$

Many clustering algorithm assume such distances given and set about constructing clusters of objects within which the distances are small. The choice of distance function is no less important than the choice of variables to be used in the study. A serious difficulty in choosing a distance lies in the fact that a clustering structure is more primitive than a distance function and that knowledge of clusters changes the choice of distance function. Thus a variable that distinguishes well between two established clusters should be given more weight in computing distances than a "junk" variable that distinguishes badly.

## 2. Euclidean Distance

The Euclidean distance between the I-th and K-th observations of a data matrix X is defined as

$$D(I,K) = \left[ \sum_{1 \leq J \leq p} \{X(I,J) - X(K,J)\}^2 \right]^{1/2} \tag{32}$$

where J is J-th variable. In one, two, or three dimensional space, this is just a "straight line" distance between the vectors corresponding to the I-th and K-th observations. When the variables are measured in different units, it is necessary to prescale the variabes to make their values comparable or, equivalently, to compute a weighted Euclidean distance.

$$D(I,K) = \left[ \sum_{1 \leq J \leq p} W(J)(X(I,J) - (X(K,J))^2 \right]^{1/2} \quad (33)$$

This form of distance is not necessary if all variables are measured on the same scale. However, even in this case, weights might be used to increase or decrease the importance of same variable. Various weighting schemes have been utilized in practice. One common weighting scheme lets $W(J)$ be the reciprocal of the variance of variable $J$ .

A general class of squared distance functions is provided by utilizing positive definite quadratic forms. Specifically, if $p$ represents a p-dimensional observation to be assigned to one of $s$ groups, then to measure the squared distance between the observation $\beta$ and the centroid (mean vector) of the $i$-th group one may consider the function

$$D_i = (\beta - \bar{x}_{i.})^T M (\beta - \bar{x}_{i.})$$

where $M$ is a positive definite matrix to ensure that $D_i \leq 0$ . Different distance functions are represented by different choices of the matrix $M$ . When $M = I$ (the identify matrix) the resulting metric is the standard Euclidean distance. Distances with the Euclidean metric are shown in Figure 1a. The variance within the data may make the unweighted Euclidean metric inappropriate. As shown on the Figure 1b, where $X$ has a larger variance than $Y$ ,

36

one may wish to weight a deviation in the  X  direction less than an equal deviation in the  Y  direction.  This is a weighted Euclidean distance frunction which makes point  A and  B  equidistance from the origin.  In this case, the matrix  M  is diagonal elements which are the reciprocals of the variances of the different variables.

Extending this idea further, it may be possible to consider the covariance among variables as well.  Figure 1c shows how the axis may be rotated so that the major axis is oriented in a direction of reflecting the positive correlation between  X  and  Y .  Again, points on the same ellipse are considered equidistance from the origin.  The matrix  M  in this case is the inverse of the covariance matrix.

Further extension of this concept will expalin some sort of generalized distance function.  If  $C_i$  represents the covariance matrix of the  i th  cluster then the distance function

$$D_i = (\beta - \bar{x}_{i.})^T \, C_i^{-1} \, (\beta - \bar{x}_{i.})$$

uses the appropriate covariance structure when determining the distance to a particular cluster centroid.  Since  $C_i$ changes to reflect the dispersion internal to each particular cluster, the use of this metric exploits differences in the dispersion characteristics of the different groups.  As shown on Figure 1d, not how a new observation (denoted by u) is

37

1a. Euclidean measure
of squared distance.

1b. Measure of squared distance
with different weights for variables.

1c. Generalized squared
distance measure.

1d. Classification when within-
group dispersions are different.

Figure 1. Euclidean Distance

closer to the centroid of group one (G1) in terms of Euclidean distance but is more likely to be assigned to group two (G2) when using the $C_i$ matrix.

### 3. Mahalanobis Distance

Another choice for the M matrix in equation (1) is $p^{-1}$ where P represents the pooled within groups covariance matrix of all the clusters.

$$P = \frac{1}{\sum\limits_{i=1}^{G} (n_i - 1)} W \qquad (34)$$

where

$$W = \sum\limits_{k=1}^{G} W_k$$

This distance is the well known Mahalanobis distance. Note that P does not change from group to group. To ensure the non-singularity of P it must be true that $p \leq (N - G)$, where N represents the total number of observations over all groups. Rewriting the distance,

$$D_i = (\beta - \bar{x}_{i.})^T W^{-1} (\beta - \bar{x}_{i.}) \qquad (35)$$

defines a distance between mean vectors $\beta$ and $\bar{x}_{i.}$ and common covariance matrix W . The Mahalanobis distance function adjusts for both scale of measurement of the variables and covariation among the variables. Use of this

39

metric is equivalent to computing distances on variables transformed to their principal components. This metric is invariant under any nonsingular transformation of original variables. For consider the transformation

$$Y = BX \qquad\qquad (36)$$

and let $D(Y_i, Y_j)$ represent Mahalanobis distance between $Y_i$ and $Y_j$.

$$
\begin{aligned}
D(Y_i, Y_j) &= (Y_i - Y_j)^T P_Y^{-1} (Y_i - Y_j) \\
&= (BX_i - BX_j)^T P_{\bar{Y}}^{-1} (BX_i - BX_j) \\
&= (X_i - X_j)^T B^T P_Y^{-1} B (X_i - X_j) \\
&= (X_i - X_j)^T B^T (BP_X B^T)^{-1} B (X_i \quad X_j) \\
&= (X_i - X_j)^T P_X^{-1} (X_i - X_j) \\
&= D(X_i, X_j)
\end{aligned}
$$

Some other common metrics are listed below:

1.  $L_1$ norm (City Block)

$$D(X_i, X_j) = \sum_{k=1}^{P} |X_{ki} - X_{kj}|$$

2.  $L_p$ norm (Minkowsky Metrics)

$$D(X_i, X_j) = \left( \sum_{k=1}^{P} |X_{ki} - X_{kj}|^P \right)^{1/p}$$

3. Uniform norm

$$D(X_i, X_j) = \underset{k = 1, 2, \ldots, p}{\text{Superemum}} \{|X_{ki} - X_{kj}|\}$$

## C. HIERARCHICAL CLUSTERING

### 1. General

The previously discussed distance measures may be used to construct a similarity matrix describing the length of all pairwise relationships among the entities (variables or data units) in the data set. The methods of hierachical cluster analysis operate on this similarity matrix to construct a tree depicting specified relationships among the entities. As shown on Figure 2, the branches on the left each represent one entity while the root represents the entire collection of entities. Moving down the tree from the branches toward the root depicts increasing aggregation of the entities into clusters. Hierarchical clustering methods which build a tree from branches to root often are called agglomerative methods.

Once a tree is constructed for N entities, the analyst may choose from as many as N sets of clusters. These clusters are nested. From the agglomerative view, when two entities are merged they are joined togehter permanently and considered as one entity for later merges; from the divisive view, when a group of entities is split into two parts, the parts are separated permenently and may be treated independently for the remainder of the analysis.

41

Figure 2. Tree for Hierarchical Clustering

Herein lie both the strength and weakness of
hierarchical methods: by taking early decisions as perma-
nent, the number of posibilities that need be examined is
reduced greatly as compared with complete enumeration;
but this same convention precludes discovering early
mistakes or capitalizing on later opportunities.

There are three major hierarchical clustering concepts:

1. Linkage Methods

2. Centroid Methods

3. Error sum of squares or variance methods.

All of these methods are suitable for clustering data units.
However, only the linkage methods are considered in this
research.

2. The General Agglomerative Procedure

Let $s_{ij}$ be the similarity between entities i
and j as defined by one of the distance measures previously
discussed. Assuming that the similarity is symmetric, the
complete schedule of similarities for all $\binom{N}{2} = \frac{1}{2}N(N - 1)$

possible pairwise combinations of entities may be arrayed
in a lower triangular similarity matrix as in Figure 3.
The $s_{ij}$ entries are nonnegative. This limitation is of
consequence only for correlation and the cosine of the angle
between vectors; the distinction between positive and nega-
tive association cannot be utilized in these clustering
methods.

$$
S = \begin{array}{|llll}
s_{21} & & & \\
s_{31} & s_{32} & & \\
s_{41} & s_{42} & s_{43} & \\
\cdot & \cdot & \cdot & \\
\cdot & \cdot & \cdot & \cdot \\
\cdot & \cdot & \cdot & \cdot \\
s_{n1} & s_{n2} & s_{n3} & \cdots\ s_{n(n-1)}
\end{array}
$$

Figure 3.  Lower Triangular Similarity Matrix

A simple remedy is to use the absolute value or the square of
the measure if it can assume negative values.  Once the
matrix is defined, the process of clustering entities is
almost trivially simple.  The general procedure for agglomer-
ative clustering on a data matrix is as follows:

(1)  Begin with  n  clusters each consisting of
exactly one entity.  Let the clusters are
labled with the numbers 1 through N.

(2)  Search the similarity matrix for the most
     similar pair of clusters.  Let the chosen
     clusters be labeled  p  and  q  and let
     their associated similarity be  $s_{pq}$, $p > q$.

(3)  Reduce the number of clusters by 1 thorugh
     merger of clusters  p  and  q .  Label the
     product of the merger  q  and update the
     similarity matrix entities in order to
     reflect the revised similarities between
     cluster  q  and all other existing clusters.
     Delete the row and column of  S  pertaining
     to cluster  p .

(4)  Perform steps 2 and 3 a total of N-1 times
     (at which point all entities will be one
     cluster).  At each step record the identity
     of the clusters which are merged and the value
     of similarity between them in order to have
     a complete record of the results.

Different agglomerative methods are implemented by
varying the procedures used for defining the most similar
pair at step 2 and for updating the revised similarity
matrix at step 3.  The similarity matrix is a given array
of numbers.  The numerical execution of the clustering
procedures is completely independent of how the similarity
values were generated or whether the entities to be
clustered are variables or data units.  However, it is
necessary to make a direct distinction between distance-like

44

measures (the smallest values correspond to the most similar pairs) and correlation-like measures (the largest values correspond to the most similar pairs); the essential difference is whether the search for the most similar pair involves seeking the minimum or maximum entry in the similarity matrix.

### 3. Single Linkage

The method of single-linkage cluster analysis is the simplest of all hierarchical techniques. At each stage, after clusters p and q have been merged, the similarity between the cluster (labeled t) and some other r is determined as follows:

1. If $s_{ij}$ is the distance-line measure

$$s_{tr} = \min (s_{pr}, s_{qr}) \qquad (37)$$

The quantity $s_{tr}$ is the distance between the two closest members of clusters t and r . If clusters t and r were to be merged, then for any entity in the resulting cluster the distance to its nearest neighbor would be at most $s_{tr}$ .

2. If $s_{ij}$ is a correlation-like measure

$$s_{tr} = \max (s_{pr}, s_{qr}) \qquad (38)$$

The quantity $s_{tr}$ is the similarity between the two most similar entities in clusters t and r . If clusters t

and  r  were to be merged, then for any entity in the resulting cluster there would be at least one other entity in the same cluster such that the pair would have a similarity at least as large as $s_{tr}$ .

The method is known as single linkage because clusters are joined at each stage by the single shortest or strongest link between them.  Since the updating process involves choosing only the minimum or maximum single-linkage clustering is invariant to any transformation which leaves the ordering of the similarities unchanged; that is, any monotonic transformation.

### 4.  Complete Linkage

The complete-linkage method is related to the single-linkage method and is no more difficult to execute.  At each stage, after clusters  p  and  q  have been merged, the similarity between the new cluster (labeled t) and some other cluster  r  is determined as follows:

1.  If  $s_{ij}$  is distance-like measure

$$s_{tr} = max\ (s_{pr}, s_{qr}) \tag{39}$$

The quantity  $s_{tr}$  is the distance between the most distant members of clusters  t  and  r .  If clusters  t  and  r  were merged, then every entity in the resulting cluster would be no farther than  $s_{tr}$  from every other entity in the cluster. The value of  $s_{tr}$  is the diameter of the samllest sphere which can enclose the cluster resulting from the merger of clusters  t  and  r.

46

2.  If $s_{ij}$ is a correlation-like measure

$$s_{tr} = \min (s_{pr}, s_{qr}) \qquad (40)$$

The quantify $s_{tr}$ is the similarity between the two most dissimilar entities in clusters  t  and  r .  If clusters t  and  r  were to be merged, then every entity in the resulting cluster would have a similarity of at least $s_{tr}$ with every other entity in the cluster.

The method is called complete linkage because all entities in a cluster are linked to each other at some maximum distance or minimum similarity.  Such a cluster is called a "maximally connected subgraph" in graph theory. In contrast to the single-linkage method, interpretation of the clusters can be made only in terms of the relation- ships within individual clusters; there is no particularly useful interpretation involving the differences between clusters.  Like the single-linkage method, complete-linkage cluster analysis is invariant to monotonic transformations of the similarity measure.  Johnson (1967) discusses this property in both single and complete linkage methods.

D.  NONHIERARCHICAL CLUSTERING

Nonhierarchical clustering methods are designed to cluster data units into a single classification of  g clusters, where  g  either is specified a priori or is determined as a part of the clustering method.  The central idea in most of these methods is to choose some initial

47

partition of the data units and then alter cluster member-
ships so as to obtain a better partition. The various
algorighms which have been proposed differ as to what
constitutes a "better partition" and what methods may be used
for achieving improvements.

The broad concept for these methods is very similar
to that underlying the steepest descent algorithms used
for unconstrained optimization in nonlinear programming.
Such algorithms begin with an initial point and then
converge to a local optimum, moving one step at a time,
the value of the objective function improving at each step.

The methods of nonhierarchical clustering typically
may be used with much larger problems than the hierarchical
methods because it is not necessary to calculate and store
the similarity matrix; it is not even necessary to store
the data set. In general, the data units are processed
serially and can be read from tape or disk as needed. This
characteristic makes it possible, at least in principle, to
cluster arbitrary large collections of data units.

In this research, we consider only the partitioning
method known as "K-MEANS" which was developed by MacQueen
(15). He used the term "K-MEANS" to denote the process of
assigning each data unit to that cluster (of k clusters)
with the nearest centroid (mean vector). The cluster
centroid changes with each transfer of an observation.

The decomposition of the total scatter matrix into
within and between groups matrices suggests possible

optimality criteria to be used in a clustering algorithm. One would like the within-groups scatter to be small relative to the between-groups scatter. Various trial clusterings could be formed using the W and B matrices as a basis for the optimality criteria which determine the best clustering. A possible choice for a criterion is to minimize trace W over all partitions into g groups. Since T is constant over all partitions, minimizing trace W is equvalent to maximizing traces B since

$$\text{trace } T = \text{trace } W + \text{trace } B \qquad (41)$$

Although trace W is invariant under an orthogonal transformation, it is not invariant under other non-singular linear transformations.

McRae (16) points out that trace W equals the total within group sum of squares, hence the "minimum variance partition" cluster solution is found by minimizing trace W .

Considerable study has been developed to alternative criteria such as those based on multivariate statistical analysis techniques, especially the methods of linear discriminant analysis and multivariate analysis of variance. Assuming the p variables are not linearly dependent, then as long as $p = N - g$ , W is positive definite symmetric and so is $W^{-1}$ . Attempts to make B and W as different as possible lead one to solving the determinantal equation:

$$|B - \lambda W| = 0 \qquad (42)$$

49

The solutions $\lambda_i$ are the eigenvalues of the matrix $W^{-1}B$ as in discriminant analysis. There are  t  non-zero eigenvalues, where  t  is the minimum of  p  and  g-1 . This is a consequence of the fact that, if  g  is less than p , the  g  group means are considered in a (g-1)-dimensional hyperplane. When  g = 2  the analysis is equivalent to two-group discriminant analysis. Linear discriminant analysis would take the vectors originally described in p-dimensional coordinate system and transform the basis to a t-dimensional system. Maximizing the largest of these eigenvalues is a criterion suggested by S.N. Roy and maximizing the trace of  $W^{-1}B$ , however is a criterion suggested by Hotelling. In both cases, large values for these statistics are sought in clustering algorithms since large values indicate large differences among (between) groups. Minimizing the ratio of determinants  $|W| \div |T|$  is a criterion widely known as Wilks' lambda discussed in the discriminant analysis. Since  T  is the same for all partitions, this criterion is equivalent to minimizing determinant  W . Both trace  $W^{-1}B$  and  $|T| \div |W|$   may be expressed in terms of the eignevalues of  $W^{-1}B$ .

$$\left|\frac{T}{W}\right| = \prod_{i=1}^{t} (1 + \lambda_i) \tag{43}$$

$$\text{trace} \quad W^{-1}B = \sum_{i=1}^{t} \lambda_i \tag{44}$$

where $t = \min(p, g-1)$. Therefore minimizing det $W$ is equivalent to maximizing $\pi(1 + \lambda_i)$.

Friedman and Rubin (6) describe the advantages of the various criteria. Those based on multivariate statistical considerations (all but trace $W$) are invariant under changes in scale for varibles (non-singular linear transformation). In fact, they are the only invariants for $W$ and $B$ under such transformations. In addition, the multivariate criteria may take into account covariation among the variables.

## V. ANALYSIS OF MULTIVARIATE UTILITY DATA

To illustrate hierarchical clustering we applied the technique described in the previous chapter to partition a set of twenty six attributes of a close-air support weapon system into a smaller collection of "superattributes". As part of an effort to evaluate the military utility of a proposed alternative U.S. Marine Corps air support rada system, AN-TPQ/27. Barr and Richards (4) extracted 26 attributes of the TPQ-27 and a baseline system, the AN-TPQ/10, and then had members of the Operational Test and Evaluation Team assess the utility of the TPQ/27 relative to that of the TPQ/10. In order that the additive model used to combine unidimensional relative utilities into a system relative utility be justifiable, it is necessary that the utilities satisfy certain independence properties described in Keeney and Raiffa (12).

Because those independence properties are very difficult for decision makers to verify for complex alternatives like the weapon systems under study, Professors Barr and Richards attempted instead to work with the attributes to try to generate a new collection which would likely satisfy, at least approximately, the conditions required to justify the additive model.

The original collection of 26 attributes is as follows:

1.  Portability

2.  Durability

3.  Time to Set Up

4.  Time to Take Down

5.  Ease of Assigning Aircraft to Targets

6.  Number of Aircraft Controlled

7.  Number of Targets

8.  Communications

9.  Mission Flexibility

10. ASRT Survivability

11. Time to Locate and Acquire Aircraft

12. Accuracy of Tracking

13. Accuracy of Delivery

14. Range

15. Aircraft Vulnerability

16. Aircraft Attack Throughout

17. Base of Adjustment and Evaluation of Results

18. Accuracy of Feedback

19. Ease of Operation

20. Man-Machine Compatibility

21. Training Requirements

22. Reliability

23. Maintainability

24. Supportability

25. Availability

26. Documentation

where $a_i$ represents an attribute $i$ and

$$I(x_{ij}, x_{kj}) = \begin{cases} 1 & \text{if } x_{ij} = x_{kj} \\ 0 & \text{if } x_{ij} \neq x_{kj} \end{cases}$$

It is easy to verify that D is a metric as defined in Chapter IV. Since we will actually work with a similarity measure in the hierarchical cluster procedure, we define the similarity between two attributes $a_i$ and $a_k$ as

$$S(a_i, a_k) = \sum_{j=1}^{12} I(x_{ij}, x_{kj}) \qquad (46)$$

One can see from this definition that the similarity between two attributes $a_i$ and $a_k$ is simply the number of team members who placed attributes $a_i$ and $a_k$ in the same partition. For example,

$$S(a_1, a_2) = 0 + 1 + 0 + 1 + 0 + 0 + 0 + 1 + 0 + 0 + 0 + 1 + 0 = 4$$

Either S or D can be used in the computer program shown in Appendix A for hierarchical clustering. One need only indicate whether he wants a correlation-like (larger values imply more similar) measure or a distance-like measure (smaller values imply more similar). We selected to use the former method. The similarity matrix extracted from

54

Table II. Data Matrix

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 6 | 2 | 1 | 3 | 2 | 4 | 3 | 1 | 3 | 1 | 1 | 1 |
| 2 | 1 | 2 | 2 | 3 | 1 | 3 | 2 | 1 | 1 | 5 | 1 | 2 |
| 3 | 6 | 1 | 1 | 5 | 2 | 4 | 3 | 1 | 6 | 1 | 2 | 1 |
| 4 | 6 | 1 | 1 | 5 | 2 | 5 | 3 | 1 | 6 | 1 | 2 | 1 |
| 5 | 2 | 5 | 3 | 1 | 3 | 1 | 1 | 2 | 2 | 2 | 3 | 3 |
| 6 | 2 | 7 | 4 | 1 | 4 | 1 | 1 | 2 | 4 | 2 | 3 | 3 |
| 7 | 2 | 7 | 4 | 1 | 4 | 1 | 1 | 2 | 4 | 2 | 3 | 3 |
| 8 | 3 | 5 | 4 | 7 | 5 | 6 | 1 | 3 | 1 | 3 | 3 | 6 |
| 9 | 2 | 5 | 4 | 1 | 3 | 1 | 1 | 2 | 4 | 2 | 3 | 3 |
| 10 | 9 | 6 | 5 | 6 | 5 | 7 | 8 | 3 | 7 | 3 | 2 | 4 |
| 11 | 2 | 8 | 6 | 4 | 3 | 2 | 1 | 4 | 4 | 2 | 3 | 3 |
| 12 | 3 | 8 | 7 | 4 | 4 | 2 | 6 | 5 | 4 | 6 | 4 | 3 |
| 13 | 3 | 8 | 7 | 4 | 4 | 2 | 6 | 5 | 4 | 6 | 4 | 3 |
| 14 | 3 | 8 | 4 | 4 | 4 | 2 | 6 | 5 | 4 | 2 | 4 | 3 |
| 15 | 7 | 6 | 5 | 6 | 5 | 7 | 7 | 3 | 7 | 3 | 7 | 4 |
| 16 | 8 | 8 | 6 | 1 | 4 | 1 | 7 | 4 | 4 | 2 | 7 | 3 |
| 17 | 8 | 8 | 8 | 4 | 3 | 2 | 5 | 6 | 5 | 4 | 5 | 3 |
| 18 | 4 | 8 | 8 | 4 | 6 | 2 | 5 | 6 | 5 | 6 | 5 | 3 |
| 19 | 4 | 5 | 3 | 1 | 3 | 1 | 1 | 2 | 2 | 4 | 3 | 3 |
| 20 | 4 | 5 | 3 | 3 | 3 | 3 | 3 | 1 | 2 | 4 | 8 | 3 |
| 21 | 5 | 4 | 9 | 2 | 3 | 8 | 4 | 7 | 2 | 4 | 6 | 5 |
| 22 | 1 | 3 | 2 | 3 | 1 | 3 | 2 | 7 | 1 | 5 | 6 | 2 |
| 23 | 1 | 3 | 9 | 3 | 1 | 3 | 2 | 7 | 1 | 5 | 6 | 2 |
| 24 | 1 | 3 | 4 | 3 | 1 | 3 | 2 | 7 | 3 | 5 | 6 | 2 |
| 25 | 1 | 3 | 2 | 3 | 1 | 3 | 2 | 7 | 1 | 5 | 6 | 2 |
| 26 | 5 | 4 | 9 | 2 | 6 | 8 | 4 | 7 | 2 | 4 | 6 | 5 |

from the data is shown in Table V-3. We present only lower triangular elements since $S(a_i, a_i) = 12$ for all $i$ and the matrix is symmetric; i.e., $S(a_i, a_k) = S(a_k, a_i)$. Zero values are not written.

The results from the hierarchical clustering are shown in Figure 4. The numbers printed along the left hand margin refer to the attribute numbers. As you proceed to the right through the tree you will observe numbers greater than 26. These correspond to the clusterings that takes place from one step to the next. For example, the number 27 shown at the juncture of 25 and 22 means that the first attribute clustered together should be 25 and 22 (this is the most similar pair). This combination is then considered as a new attribute which is later combined with the attribute 30 (itself a combination of 23 and 24) to form the attribute 31. This is later combined with attribute 2 to form attribute 40, etc.

As discussed in Chapter IV a decision has to be made as to how many clusters (superattributes) are desired. All hierarchical methods will continue clustering until there is a single cluster. In order to decide on the number of clusters (and their composition) one need only image drawing a vertical line through the tree at various places. Each intersection of the tree with the vertical line results in a cluster. For example, teh vertical line at the point A results in the 6 clusters shown in Table V-4.

It is clear from observing the above collection that some of the attributes are highly correlated and nonredundant. If one tries to assign an importance weights to each attributes separately, there is a distinct likelihood that some of the overlapping strongly into related attributes might effectively be double or triple weighted or more producing biased result. It is an effort to prevent this from happening, Barr and Richards aksed the utility assessment team to partition the 26 attributes into a smaller collection in such a way that attributes within a group are similiar and attributes in different groups are unrelated the sense that utility assessments for attributes in one group do not depend on the amounts of attributes in any other group.

The total number of groups was not prespecified. Instead, each team member was allowed to partition the 26 attributes into any number of groups. The resulting multivariate data array is shown in Table V-2. An element $x_{ij}$ is the number of the group into each team member $j$ put attribute $i$.

Let us define a distance measure for this data array as follows:

$$D(a_i, a_k) = \sum_{j=1}^{12} (1 - I(x_{ij}, x_{kj}) \qquad (45)$$

Table III.  Similarity Matrix for Superattribute Determination.

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 2 | 4 | | | | | | | | | | | | | | | | | | | | | | | | | |
| 3 | 8 | 1 | | | | | | | | | | | | | | | | | | | | | | | | |
| 4 | 7 | 1 | 11 | | | | | | | | | | | | | | | | | | | | | | | |
| 5 | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 6 | | | | | 8 | | | | | | | | | | | | | | | | | | | | | |
| 7 | | | | | 8 | 12 | | | | | | | | | | | | | | | | | | | | |
| 8 | | 1 | | | 3 | 3 | 3 | | | | | | | | | | | | | | | | | | | |
| 9 | | | | | 10 | 10 | 10 | 3 | | | | | | | | | | | | | | | | | | |
| 10 | | | 1 | 1 | | | | 3 | | | | | | | | | | | | | | | | | | |
| 11 | | | | | 6 | 6 | 6 | 2 | 7 | | | | | | | | | | | | | | | | | |
| 12 | | | | | 1 | 3 | 3 | 3 | 1 | 2 | | 5 | | | | | | | | | | | | | | |
| 13 | | | | | 1 | 3 | 3 | 3 | 1 | 2 | | 5 | 12 | | | | | | | | | | | | | |
| 14 | | | | | 2 | 5 | 5 | 2 | 4 | | | 5 | 9 | 9 | | | | | | | | | | | | |
| 15 | | | | | | | | 3 | | 9 | | | | | | | | | | | | | | | | |
| 16 | | | | | 3 | 5 | 5 | | 4 | 1 | 6 | 3 | 3 | 4 | 3 | | | | | | | | | | | |
| 17 | | | | | 2 | 1 | 1 | 1 | 2 | | 5 | 4 | 4 | 3 | | 2 | | | | | | | | | | |
| 18 | | | | | 1 | 1 | 1 | 1 | 1 | | 5 | 4 | 4 | 3 | | 1 | 9 | | | | | | | | | |
| 19 | | | | | 9 | 5 | 5 | 2 | 7 | | 3 | | | | | 2 | 2 | 1 | | | | | | | | |
| 20 | 2 | 3 | 1 | 1 | 5 | 1 | 1 | 1 | 3 | | 2 | | | | | | 2 | 1 | 8 | | | | | | | |
| 21 | | | | | 1 | | | | 1 | | 1 | | | | | | 2 | | 3 | 3 | | | | | | |
| 22 | 2 | 8 | | | | | | 1 | | | | | | | | | | | | 2 | 2 | | | | | |
| 23 | 2 | 7 | | | | | | 1 | | | | | | | | | | | | 2 | 3 | 11 | | | | |
| 24 | 3 | 6 | | | | | | 1 | | | | | | . | | | | | | 2 | 3 | 10 | 11 | | | |
| 25 | 2 | 1 | | | | | | | | | | | | | | | | | | 2 | 2 | 12 | 11 | 10 | | |
| 26 | | | | | 1 | | | | | | | | | | | | 1 | 1 | 2 | 2 | 9 | 3 | 4 | 4 | 3 | |

Figure 4. Tree for 26 Attributes

59

The superattributes used in the utility study are those sohwn in Table IV.  A careful examination of the attributes which comapre the clusters shows that the results so obtained are intuitively agreeable.  The names supplied to the superattribures are somewhat natural descriptions of the clusters obtained.

The listing of the computer program and sample output are given in Appendix A.

Table IV.    Superattributes

| Superattributes | | Component Attributes |
|---|---|---|
| Facility of movement | 1. | Portability |
| | 3. | Time to set up |
| | 4. | Time to take down |
| Facility of Use | 5. | Ease of assigning aircraft to targets |
| (precision) | 6. | Number of aircraft controlled |
| | 7. | Number of targets |
| | 9. | Mission flexibility |
| | 11. | Time to locate and acquire aircraft |
| | 16. | Aircraft attack throughput |
| | 17. | Ease of adjustment |
| | 18. | Accuracy of feedback |
| | 19. | Ease of operation |
| | 20. | Man-machine compatibility |
| | 12. | Accuracy of tracking |
| | 13. | Accuracy of delivery |
| | 14. | Range |
| Survivability | 10. | ASRT Survivability |
| | 15. | Aircraft vulnerability |
| Learning | 21. | Training requirements |
| | 26. | Documentation |
| Readiness | 2. | Durability |
| | 22. | Reliability |
| | 23. | Maintainability |
| | 24. | Supportability |
| | 25. | Availability |
| Communications | 8. | Communications |

# VI. ANALYSIS OF ARMY TANK DATA

## A. DATA STRUCTURE

In order to illustrate the nonhierarchical clustering methodology, principal components analysis, and discriminant analysis data on Army tanks from eight different countries were taken from Jane's Book of Weapon Systems (1979-80). A total of twenty-four tanks were included in the data array with observation on each of 10 variables. The 10 variables are listed below:

1. Weight (ton)

2. Length (meter)

3. Width (meter)

4. Height (meter)

5. Road Speed (kilometer per hour)

6. Trench Crossing (meter)

7. Ground Pressure ($Kg/cm^2$)

8. Maximum Armament (rounds)

9. Ground Clearance (meter)

10. Power to Engine Ratio (BHP/ton)

The twenty-four tanks and the associated countries are shown below:

| Identification Number | Type/Name | Country |
| --- | --- | --- |
| 11 | T-62 | |
| 12 | T-54 | U.S.S.R. |

| Identification Number | Type/Name | Country |
|---|---|---|
| 13 | T-10 | |
| 14 | ASU-85 | |
| 15 | MK-5/Chieftain | |
| 16 | MK-3/Vickers | |
| 17 | MK-13/Centurion | U.K. |
| 18 | CVR(T)/Scorpion | |
| 19 | XM-1 | |
| 20 | M60A2 | |
| 21 | M60 | U.S.A. |
| 22 | M48 | |
| 23 | M47 | |
| 24 | PZ61 | SWITZER-LAND |
| 25 | PZ68 | |
| 26 | STRV-103 | SWEDEN |
| 27 | Ikv-91 | |
| 28 | TYPE61 | JAPAN |
| 29 | TYPE74 | |
| 30 | Leopard 2 | |
| 31 | Leopard | W. GER-MANY |
| 32 | TAM | |
| 33 | AMX 30 | FRENCH |
| 34 | AMX 13 | |

We conjecture that a cluster analysis of the tank data will
result in clusters corresponding to nationality since the
nations may have different emphasis on the variables in
the design of their tanks.

B.  NONHIERARCHICAL CLUSTER ANALYSIS OF TANK DATA

  1.  The MIKCA Algorithm

      The specific algorithm chosen for the nonhierarchical
cluster analysis for the tank data is the MIKCA (Multivariate
Iterative K-MEANS Clustering Algorithm) program written by
Douglas J. McRae as a part of his doctoral dissertation
at the University of North Carolina, Chapel Hill.

      Reference to the flow chart in Figure 5 will aid the
reader in following discussion of the algorithm.  Inputs to
program are the data matrix, an estimate for  $g$  (the
number of clusters), and choice of criterion and distance
functions.

      In the first step, preliminary claculations are made,
such as the variable means  and standard deviations, as
well as the cross product matrix  $T$ .  The next step forms
the initial cluster centers.  Then each of the other
observations is assigned to the nearest cluster.  Euclidean
distance is used for this initial phase, and the cluster
centroids are recomputed after each observation is assigned
to a group.  The observations are considered in the same
order as they were input.  After all of them have been
assigned to clusters, the criterion value is computed.
This initial cluster-finding technique is referred to as a

Figure 5. MIKCA Flow Chart

one-pass K-MEANS procedure. It is performed three times, and the solution which yields the best criterion value is chosen as the initial cluster solution.

After the initial solution has been found, the program advances to the iterative K-MEANS phase where the observations are again considered in the order in which they were input to the program. It is this phase where the user's choice of distance function is used. The distance from each observation to each cluster centroid is again computed, this time with the user's distance function, the assignment to the closest centroid being made and the centroid updated to reflect its new membership. After considering all n observations in this manner, the new criterion value is checked for possible improvement during the K-MEANS iteration. As long as the criterion value improves, the K-MEANS procedure is repeated; if the criterion fails to improve then the MIKCA algorithm goes to the next step, the individual switches section. Note the importance of the order of consideration of the observations. The order is important because the cluster means are recomputed after each observation is reassigned.

In the individual switches phase, consideration is given to moving each observation to every other cluster, the move being made if and only if an improvement in the value of the criterion results. An elaborate labelling procedure provides a unique order in which to consider each observation.

This procedure continues until a complete pass through the data is made with no changes in cluster membership.

The MIKCA alogorithm provides the following options for distance and criterion functions.

Criterion

1. Minimum trace W

2. Minimum determinant W

3. Maximum largest order of $|B - \lambda W| = 0$

4. Maximum sum of roots of $|B - \lambda W| = 0$

Distance

1. Euclidean

2. Weighted Euclidean

3. Mahalanobis

A complete computer program is listed in Appendix B.

2. Cluster Results for Tank Data

For clustering of the tank data we selected the minimum trace W criterion and the weighted Euclidean distance function. The algorithm automatically provides weights for the weighted Euclidean distance function. The results of the clustering with four clusters are shown in Table V.

The conjecture of clustering by nationalities is supported by the results. The three Soviet tanks make up one cluster and the two British and four of the United States tanks were found to be similar. A third cluster consists of four tanks which are very lightweight. The

Table V. THE FINAL CLUSTER SOLUTION IS

```
CLUSTER 1
  SIZE = 11
  CENTER =   40.60      7.00      3.26  53.44      2.55  59.18
                        2.64      0.66            0.46  20.70

  THE OBSERVATIONS ARE
   16   19   24   25
   28   29   30
   33

CLUSTER 2
  SIZE = 6
  CENTER =   51.87      7.15      3.55  59.50      3.06  43.05
                        2.80      0.66            0.43  15.55

  THE OBSERVATIONS ARE
   15   17   20   21
   23

CLUSTER 3
  SIZE = 4
  CENTER =   13.00      5.40      2.77  45.75      2.21  64.27
                        2.14      0.51            0.38  20.07

  THE OBSERVATIONS ARE
   14   18   27   34

CLUSTER 4
  SIZE = 3
  CENTER =   41.00      6.85      3.50  27.67      2.35  44.33
                        2.65      0.77            0.44  15.07

  THE OBSERVATIONS ARE
   11   12   13
```

final cluster consists of the rest of the tanks, including tanks of United States allies from West Germany, France, Sweden, Switzerland and Japan.

A natural question to ask after observing the results of a cluster analysis is what variables most strongly influence the clustering that was observed. A clue is provided by the composition of the cluster containing all of the lightweight tanks. This suggests that weight is an important distinguishing feature. This is examined in the principal components analysis and the discriminant analysis in the next two section.

C.  PRINCIPAL COMPONENTS ANALYSIS

The Statistical Package for Social Sciences (SPSS) (14) subprogram FACTOR was used for the principal components analysis. It is designed both for the factor analysis and the principal components analysis. The outputs are designed to be self-explanatory. In this example, the first 5 components accoung for 90% of the variance and the remaining components account for only 10% of the variance (Figure 6).

The subprogram FACTOR provides a graphical presentation (Figure 7) for the factors that have been determined by the orthogonal rotations (in this example, variance maximization rotation). In reading the graphs, one should be attentive to following three features: (1) the relative distance of a variable from the axis, (2) the direction of a variable in relation to the axis, and finally (3) clustering of variables and their relative position to each other.

PRINCIPAL COMPONENT ANALYSIS OF TANKS

FILE NONAME (CREATION DATE = 03/19/80)

03/19/80

| VARIABLE | EST COMMUNALITY |
|---|---|
| A1 | 1.00000 |
| A2 | 1.00000 |
| A3 | 1.00000 |
| A4 | 1.00000 |
| A5 | 1.00000 |
| A6 | 1.00000 |
| A7 | 1.00000 |
| A8 | 1.00000 |
| A9 | 1.00000 |
| A10 | 1.00000 |

| FACTOR | EIGENVALUE | PCT OF VAR | CUM PCT |
|---|---|---|---|
| 1 | 4.78856 | 47.9 | 47.9 |
| 2 | 1.71316 | 17.1 | 65.0 |
| 3 | 1.26768 | 12.7 | 77.7 |
| 4 | 0.74450 | 7.4 | 85.1 |
| 5 | 0.55505 | 5.6 | 90.7 |
| 6 | 0.34472 | 3.4 | 94.1 |
| 7 | 0.26042 | 2.6 | 96.7 |
| 8 | 0.15890 | 2.0 | 98.7 |
| 9 | 0.08262 | 0.8 | 99.6 |
| 10 | 0.04438 | 0.4 | 100.0 |

Figure 6. Summary Table of Principal Components Analysis on Tank

70

PRINCIPAL COMPONENT ANALYSIS OF TANKS

FILE ACADE (CREATION DATE = 03/15/80)

HORIZONTAL FACTOR 1     VERTICAL FACTOR 2

03/15/80

VERTICAL FACTOR 2

10

5

3

2

1

Figure 7. Graphical Presentation

PRINCIPAL COMPONENT ANALYSIS OF TANKS

FILE ACAJPS (CREATION DATE = 03/15/80)

FACTOR SCORE COEFFICIENTS

| | FACTOR 1 | FACTOR 2 | FACTOR 3 | FACTOR 4 | FACTOR 5 | FACTOR 6 | FACTOR 7 | FACTOR 8 | FACTOR 9 | FACTOR 10 |
|---|---|---|---|---|---|---|---|---|---|---|

Figure 8. Factor Score Coefficients

For example, variables 5 (road speed) and 10 (power to engine ratio) contribute heavily to the first principal component while variables 1 (wieght) and 3 (width) contributes most strongly to the second principal component. Variables 2, 4, 6, 7, 8, 9 are not as important. The weights accorded each variable in the 10 factors (principal components) are shown in Figure 8. The complete SPSS program is listed in Appendix C.

## D. DISCRIMINANT ANALYSIS

The SPSS subprogram DISCRIMINANT was used to determine that function or those functions of the 10 variables that best discriminant among the four clusters determined in previous section.

The maximum number of discriminant functions to be derived is either one less than the number of groups or equal to the number of discriminating variables. This subprogram provides two measures for judging the importance of discriminant functions. One of these is the relative percentage of the eigenvalue associated with the function. It is a measure of the relative importance of the function. The sum of the eigenvalues is a measure of total variance existing in the discriminating variables. Since discriminant functions are derived in order of their importance, this process can be stopped whenever the relative percentage is judged to be too small. Of course, there is no fixed rule for deciding whatis too small. In this research, we selected arbitrary, a significance level of 0.10. The output shown

73

in Figure 9 suggests that we therefore consider only the first two discriminant functions.

The second measure judging the importance of a discriminant function is its associated canonical correlation. The canonical correlation is a measure of association between the single discriminant function and the set of (g-1) dummy variables which define the g group memberships. It tells us how closely the function and the group variable are related, which is just another measure of the function's ability to discriminate among the groups. From Figure 10, the first two discriminant functions are each highly corre-lated with the groups but the third has only a moderate correlation.

The next criterion for eliminating discriminant functions is to test for the statistical significance of discriminating information not already accounted for by the earlier functions. As each function is derived, starting with no (zero) functions, Wilks' lambda is computed. Lambda is an inverse measure of the discriminating power in original variables which has not yet been removed by the discriminant functions - the larger lambda is, the less is the information remaining. Lambda can be transformed into a chi-square statistic for an easy test of statistical significance. In Figure 9, Wilks' lambda was .594 after the first two functions had been derived. This corresponds to a chi-square of 8.8476 with a probability level of .1823.

74

Figure 9. Summary Table of Discriminant Analysis on Tank

TASK DISCRIMINANT

STANDARDIZED CANONICAL DISCRIMINANT FUNCTION COEFFICIENTS

|  | FUNC 1 | FUNC 2 | FUNC 3 |
|---|---|---|---|

POOLED WITHIN-GROUPS CORRELATIONS BETWEEN CANONICAL DISCRIMINANT FUNCTIONS AND DISCRIMINATING VARIABLES
(VARIABLES ARE ORDERED BY THE FUNCTION WITH LARGEST CORRELATION AND THE MAGNITUDE OF THAT CORRELATION.)

|  | FUNC 1 | FUNC 2 | FUNC 3 |
|---|---|---|---|

UNSTANDARDIZED CANONICAL DISCRIMINANT FUNCTION COEFFICIENTS

|  | FUNC 1 | FUNC 2 | FUNC 3 |
|---|---|---|---|

Figure 10.  Canonical Discriminant Function Coefficients

76

This means that a lambda of this magnitude or smaller has a
.1823 probability of occurring due to the chances of
sampling even if there was no further information to be
accounted for by a third function in the population.
Clearly, a third function is not statistically significant
in this case.

The standarized discriminant function coefficients
corresponding to the values of the $v_{ij}$'s discussed in the
previous section are used to compute the discriminant score
for a case (observation) in which the original discriminating
variables are in standard form. The discriminant score is
computed by multiplying each discriminating variable by its
corresponding coefficient and adding together these products.
There is a separate score for each observation on each
function. The coefficients have been derived in such a way
that the discriminant scores produced are in standard form.

When the sign if ignored, each standard discriminant
function coefficient represents the relative contribution of
its associated variable to that function. The sign merely
denotes whether the variable is making a positive or negative
contribution.

A graphical presentation is shown in Figure 11 using
the first and the second canonical discriminant function as
the axis. From this scatterplot, we can easily see that
Soviet tanks (labelled 1) are well distinguished from the
all of the others using only the first two discriminant

Figure 11.  Scatter Plot for all Groups

functions. Also, all the lightweight tanks are clearly separated from the others. The distinction between groups 2 and 3 is also clear though not separated from each other as much as from groups 1 and 4. The complete SPSS program for the discriminant analysis is listed in Appendix D.

# VII. CONCLUSION

The multivariate analysis techniques of cluster analysis, principal components analysis and discriminant analysis are useful in real world problems for examining observations on each of several dimension. Each of the techniques is related mathematically to the others, and each complements the other in explaining the data.

Computer software is readily available in many sources. The software used in this thesis for hierarchical clustering, principal components analysis, and discriminant analysis was from the IMSL package and SPSS. For nonhierarchical clustering, we used the FORTRAN program developed by McRae (16). All of this softw.re is readily available and documented at the Naval Postgraduate School.

Appendix A:

# HIERACHICAL CLUSTERING

THIS IS A PROGRAM FOR HIERACHOCAL CLUSTERING

```
VARIABLE DESCRIPTION
NO       INPUT NUMBER OF DATA POINTS TO BE CLUSTERED.
         ND-1 CLUSTERS ARE FORMED NUMBERED CONSECUTIVELY
         FROM ND+1 TO ND+(ND-1). ND MUST BE GREATER THAN 2.
         ND IN THE RANGE 2,3,...,500 IS ALLOWED
         (SEE REMARKS)
IOPT     INPUT OPTIONS VECTOR OF LENGTH 2
         IOPT(1)=0 IMPLIES SINGLE-LINKAGE DESIRED.
         OTHERWISE COMPLETE-LINKAGE IS DESIRED.
         IOPT(2)=0 IMPLIES THE SIMILARITIES ARE DISTANCE-
         LIKE (I.E. SMALLER IMPLIES CLOSER). POSITIVE
         THE SIMILARITIES ARE ASSUMED CORRELATION-LIKE (SEE REMARKS)
XSIM     INPUT/OUTPUT VECTOR OF LENGTH (NC+1)*NC)/2
         CONTAINING THE SIMILARITY MATRIX IN SYMMETRIC
         STORAGE MODE. XSIM(((I-1)*I)/2+J) CONTAINS THE SIMILARITY OF
         XSIM((I-1)*I)/2+J) GREATER THAN J.
         THE I-TH AND J-TH DATA POINTS. ON INPUT THE
         DIAGONAL ELEMENTS (SIMILARITY OF ITSELF) ARE
         ARBITRARY AND NEED NOT BE DEFINED.
         XSIM IS DESTROYED ON OUTPUT.
CLEVEL   OUTPUT VECTOR OF LENGTH ND-1. CLEVEL(K) CONTAINS THE
         SIMILARITY LEVEL AT WHICH CLUSTER NC+K WAS FORMED.

DIMENSION XSIM(500),ICPT(2),ICUT(2),CLEVEL(500),TITLE(20),
1ICRSCN(500),IPTR(ICC),ICLSCN(500),CLVEK(ICO),ACLFST(ICC),
2LEFTRT(100),STARST(100),YSIM(2),IAC(4),

ICLSCN   OUTPUT VECTOR FOR SUBPROGRAM USTREE.
         NUMBER ND+K WAS FORMED BY MERGING CLUSTER ICLSCN(K)
         WITH CLUSTER ICRSCN(K). THE
         THIS WOULD BE THE INPUT VECTOR FOR SUBPROGRAM USTREE.
ICRSCN   OUTPUT VECTOR OF RIGHTSONS OF LENGTH NC-1. THE
         RIGHTSON OF CLUSTER ND+K IS CONTAINED IN ICRSCN(K).
         THIS WOULD BE THE INPUT VECTOR FOR SUBPROGRAM USTREE.
IPTR     WORK VECTOR OF LENGTH ND.
IER      ERROR PARAMETER. (OUTPUT)
         TERMINAL ERROR
         IER=129 INDICATES ND WAS LESS THAN 3.
INC      INPUT VECTOR OF LENGTH 4. INC(I) CONTAINS WHEN I=1,
         MUST BE BETWEEN ND AND 2*ND. EXCLUSIVELY (REMARKS)
         I=2, NUMBER OF PRINTABLE SPACES PER LINE ON THE
         OUTPUT (PRINTER) DEVICE. IAD(2) MUST EXCEED 4.
```

81

I=3, NUMBER OF HORIZONTAL SLICES OF TREE DESIRED
     TO PROVIDE THE NECESSARY DETAIL.
     INC(3) MUST BE POSITIVE.
I=4, NUMBER OF FILLER LINES PRINTED BETWEEN NODE
     LINES (USUALLY SUFFICIENT).

YSIM   INPUT VECTOR OF LENGTH 2 CONTAINING THE INTERVAL ON
       THE VERTICAL SCALE USED TO PLOT THE TREE
       (SEE VERTICAL CLEVEL. LEVEL YSIM(1) IS WHERE THE
       TERMINAL NODES (FOR IMSLCLUSTERING ROUTINES, THE
       DATA POINTS) ARE PRINTED. LEVEL YSIM(2) IS THE
       ENDPOINT YSIM(2) SHOULD INCLUDE THE LEVEL FOR THE
       HEAD NODE (INC(1)). SEE REMARKS.

NAMELIST /NAM1/XSIM
   ICUT    WORK VECTOR OF LENGTH INC(2)-4.
   CLVLSK  WORK VECTOR OF LENGTH NC.
   NCLRST  WORK VECTOR OF LENGTH NC.
   LEFTRT  WORK VECTOR OF LENGTH NC.
   STARST  WORK VECTOR OF LENGTH NC.

REMARKS 1. THE DATA CLUSTERS ARE NUMBERED 1 TO NC. THE NC-1 CLUSTERS
        FORMED BY MERGING ARE NUMBERED NC+1 TO NC+(NC-1) AND DECREASE
        IN SIMILARITY, MAKING IT EASY TO IDENTIFY THE MOST SIMILAR
        CLUSTERS.

        2. SIMILARITIES GENERALLY SHOULD BE NONNEGATIVE. RAW
        CORRELATIONS TAKE ON VALUES R WHERE R RANGES IN THE CLOSED
        INTERVAL (-1, 1), AND SIMILARITY, IF R=1 ARE
        R=-1 BOTH MEAN HIGH SIMILARITY. THEN THE TRANSFORMATION,
        RR EQUAL TO THE ABSOLUTE VALUE, CF, FR EQUAL TO R SQUARED,
        ARE APPROPRIATE. IF R=-1 REPRESENTS VERY LOW SIMILARITY THEN
        THE TRANSFORMATION, RR=1-R, BECOMES A DISTANCE-LIKE SIMILARITY.

        3. NOTE THAT THE ORIGINAL DATA MATRIX OF THE USER DOES NOT
        ENTER CCLINK, ALLOWING THE USER TO DEFINE, VIA XSIM,
        WHATEVER MEASURE OF SIMILARITY SEEMS MOST APPROPRIATE.

        4. THE TREE MAY BE TOO LARGE TO FIT IN INC(2) SPACES
        REPRESENTING THE INTERVAL (YSIM(1),YSIM(2)). IF SO, THE TREE CAN
        PRINT THE TREE IN IND(3) SLICES OF THIS INTERVAL WHICH MAY BE
        CUT AND TAPED TOGETHER.

        READ(5,11C) (TITLE(I),I=1,20)
        5. TO PRINT THE ENTIRE TREE FROM IMSL SUBROUTINE CCLINK, THE
        HEAD NODE INC(1) =NC+(ND-1).

        6. OUTPUT IS WRITTEN TO THE UNIT SPECIFIED BY IMSL ROUTINE
        UGETIO. SEE THE UGETIO DOCUMENT FOR DETAILS.

```
C     7. REVERSALS (IER=130) MAY OCCUR IN TWO WAYS. FIRST, THE ACCES
C        MAY BE JOINED AT A LOWER LEVEL (CLOSER TO YSIM(1)). SECOND, THE
C        LEVEL OF THE HEAD NODE MAY LIE ABOVE THE INTERVAL
C        (YSIM(1),YSIM(2)).
C     8. FOR PROPER DISPLAY, THE TREE CREATED BY USTREE SHOULD BE
C       -URNED TO AN UPRIGHT POSITION.
C
      WRITE(6,111) (TITLE(I),I=1,20)
C
C     READ THE NUMBER OF OBSERVATION AND
C     TWO INPUT OPTIONS.
C
      READ(5,112) ND,IOPT(1),IOPT(2)
      WRITE(6,2CC) ND
      IF(NC.LE.2) GO TO 50
   2==((ND+1)*NC)/2
      IF(IOPT(1).EQ.0) GO TO 30
      WRITE(6,113)
   25 IF(IOPT(2).EQ.0) GO TO 31
      GO TO 32
   30 WRITE(6,115)
      GO TO 25
   31 WRITE(6,116)
   32 WRITE(6,13) M
C
C     READ INPUT DATA MATRIX BY THE
C     NAME LIST FORMAT. FROM THE LOWER TRIANGULAR MATRIX.
C     AN ARRY OF VECTOR
C
      READ (5,NAM1)
      WRITE (6,300)
      JJJJ=1
      DO 10 I=1,NC
      JKK=JJJJ+I-1
      WRITE (6,101) I,(XSIM(JK),JK=JJJJ,JJKK)
      JJJJ=JJJJ+I
   10 WRITE (6,10C) (I,I=1,ND)
C
C     CALL IMSL SUBPROGRAM OCLINK
C
      CALL OCLINK(ND,ICFT,XSIM,CLEVEL,ICLSON,ICRSON,IFTF,IER)
      WRITE(6,1)
C
C     WRITE THE SIMILARITY LEVEL,ICLSON AND
C     ICRSON TO CONSTRUCT A TREE.
```

83

```fortran
      WRITE(6,2) (CLEVEL(I),ICLSON(I),ICFSON(I),I=1,25)
      IAC(1)=2*NC-1
      IAC(2)=70
      IAC(3)=1
      IAC(4)=1
      YSIM(1)=CLEVEL(1)
      YSIM(2)=0.0

C     CALL IMSL USTREE FOR THE GRAPHICAL
C     PRESENTATION.

      CALL USTREE(ND,ICLSON,ICRSON,CLEVEL,IND,YSIM,ICUT,CASES,IAC,LEFT,
     1LEFTRT,STARST,IER)
      GO TO(6,500)
6     WRITE(6,500)
500   CONTINUE
2     FORMAT(1H1,5X,'CLEVEL',6X,'ICLSON',6X,'ICRSON')
      FORMAT(2X,F10.1,2X,I7,5X,I7)
      FORMAT(/,'NUMBER OF OBSERVATIONS IN XSIM ARE',I6)
500   FORMAT(//,5X,26I4)
      FORMAT(/,I5,26F4.0)
100   FORMAT(20A4)
      FORMAT(1X,20A4)
      FORMAT(I4,2I1)
      FORMAT('0','COMPLETE-LINKAGE IS PERFORMED.')
      FORMAT('0','SIMILARITIES ARE POSITIVE CORRELATION-LIKE.')
      FORMAT('0','SINGLE LINKAGE IS PERFORMED.')
      FORMAT('0','SIMILARITIES ARE DISTANCE-LIKE.')
      FORMAT('0','NUMBER OF DATA POINTS ARE',I4)
2000  FORMAT('0','INITIAL DISTANCE MATRIX',//)
3000  FORMAT('0','YOUR DATA ARE TOO SMALL.')
      END
```

INTERACTICAL CLUSTERING ON UTILITY DATA

NUMBER OF DATA POINTS ARE 26

SINGLE LINKAGE IS PERFORMED

SIMILARITIES ARE POSITIVE CORRELATION-LIKE

NUMBER OF OBSERVATIONS IN DSIM ARE 351

INITIAL DISTANCE MATRIX

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 50. | | | | | | | | | | | | | | | | | | | | | | | | | |
| 2 | 4. | 50. | | | | | | | | | | | | | | | | | | | | | | | | |
| 3 | 8. | 1. | 50. | | | | | | | | | | | | | | | | | | | | | | | |
| 4 | 7. | 1. | 11. | 50. | | | | | | | | | | | | | | | | | | | | | | |
| 5 | 0. | 0. | 0. | 0. | 50. | | | | | | | | | | | | | | | | | | | | | |
| 6 | 0. | 0. | 0. | 0. | 8. | 50. | | | | | | | | | | | | | | | | | | | | |
| 7 | 0. | 1. | 0. | 2. | 8. | 12. | 50. | | | | | | | | | | | | | | | | | | | |
| 8 | 0. | 1. | 2. | 2. | 3. | 3. | 3. | 50. | | | | | | | | | | | | | | | | | | |
| 9 | 0. | 0. | 1. | 1. | 10. | 10. | 0. | 3. | 50. | | | | | | | | | | | | | | | | | |
| 10 | 0. | 0. | 1. | 0. | 4. | 6. | 2. | 7. | 0. | 50. | | | | | | | | | | | | | | | | |
| 11 | 0. | 0. | 0. | 0. | 1. | 3. | 1. | 2. | 0. | 5. | 50. | | | | | | | | | | | | | | | |
| 12 | 0. | 0. | 0. | 0. | 3. | 3. | 1. | 2. | 5. | 5. | 12. | 50. | | | | | | | | | | | | | | |
| 13 | 0. | 0. | 0. | 0. | 2. | 2. | 2. | 4. | 5. | 9. | 5. | | | | | | | | | | | | | | | |
| 14 | 0. | 0. | 0. | 3. | 5. | 2. | 4. | 5. | 0. | 0. | 3. | 3. | 3. | 50. | | | | | | | | | | | | |
| 15 | 0. | 0. | 0. | 0. | 0. | 1. | 1. | 0. | 4. | 3. | 3. | 0. | 0. | 50. | | | | | | | | | | | | |
| 16 | 0. | 0. | 0. | 3. | 5. | 2. | 1. | 2. | 5. | 3. | 2. | 0. | 3. | 50. | | | | | | | | | | | | |
| 17 | 0. | 0. | 0. | 2. | 1. | 1. | 0. | 4. | 2. | 2. | 0. | 2. | 1. | 50. | | | | | | | | | | | | |
| 18 | 0. | 0. | 0. | 1. | 1. | 0. | 0. | 3. | 0. | 0. | 0. | 0. | 0. | 5. | 50. | | | | | | | | | | | |
| 19 | 0. | 0. | 2. | 5. | 1. | 0. | 2. | 2. | 2. | 2. | 1. | 8. | 50. | | | | | | | | | | | | | |
| 20 | 2. | 2. | 1. | 1. | 1. | 0. | 2. | 0. | 0. | 2. | 1. | 3. | 3. | 50. | | | | | | | | | | | | |
| 21 | 0. | 0. | 0. | 0. | 0. | 1. | 0. | 0. | 0. | 0. | 0. | 2. | 3. | 3. | 50. | | | | | | | | | | | |
| 22 | 0. | 0. | 0. | 1. | 0. | 0. | 0. | 0. | 0. | 0. | 0. | 2. | 2. | 2. | 50. | | | | | | | | | | | |
| 23 | 0. | 0. | 0. | 0. | 0. | 0. | 0. | 0. | 0. | 0. | 1. | 2. | 3. | 11. | 50. | | | | | | | | | | | |
| 24 | 3. | 0. | 0. | 1. | 0. | 0. | 0. | 0. | 0. | 0. | 2. | 3. | 10. | 11. | 10. | 50. | | | | | | | | | | |
| 25 | 2. | 1. | 0. | 0. | 0. | 0. | 0. | 0. | 0. | 0. | 2. | 2. | 12. | 11. | 10. | 50. | | | | | | | | | | |
| 26 | 0. | 1. | 0. | 0. | 0. | 0. | 0. | 0. | 0. | 1. | 2. | 5. | 3. | 4. | 4. | 3. | 50. | | | | | | | | | |

ICRSON

22
12
6
23
30
3
5
21
19
17
10
14
20
1
16
44
45
40
47
46
8
50

ICLSCA

25
1
7
4
2
24
7
4
9
3
6
4
8
5
6
1
4
2
1
1
7
5
2
8
6
8

CLEVEL
...........

Appendix B:

K-MEANS ITERATIVE CLUSTERING PROGRAM BY D.J. MCRAE

THIS PROGRAM ENABLES THE USER TO CLUSTER A DATA MATRIX
UP TO SIZE (600 X20) INTO A MAXIMUM OF 20 CLUSTERS.
SEVERAL OPTIONS ARE AVAILABLE WITH RESPECT TO THE
METHOD OF COMPUTATION AND CRITERIA TO ACCOMPLISH THE
CLUSTERING. THE PROGRAM MAY BE USED ON EITHER
THE CP/CMS TERMINALS OR BY THE CARD READER. THE APPROPRIATE
METHOD FOR EACH IS DESCRIBED BELOW.

CP/CMS TERMINAL USE
THE PROGRAM SHOULD BE READ INTO CP AS IS, BUT
PRECEDED BY THE FOLLOWING CARDS, STARTING IN COLUMN ONE

    CF67USERID XXXXG
    OFFLINE READ RRREAD FORTRAN

THE DATA MATRIX CAN BE PLACED IN THE COMPUTER
EITHER BY USE OF THE OFFLINE READ OR BY TYPING IN THE
INFORMATION ON THE TERMINAL. IF OFFLINE READ IS USED THE DATA
MUST BE IN THE SPECIFIED FORMAT PRECEDED BY THE FOLLOWING
TWO CARDS, STARTING IN COLUMN ONE:

    CF67USERID XXXXG
    OFFLINE READ FILE FT04FYYY

THE PROGRAM CAN THEN BE EXECUTED BY FIRST ISSUING THE
CMS COMMAND

FOLLOWED BY       F RRREAD

                  $ RRREAD

CS/CARD READER USE

THE PROGRAM MUST CONTAIN THE FOLLOWING INFORMATION
IN THE FOLLOWING ORDER:
    STANDARD GREEN JOB CARD,TIME=(AS DESIRED)
    //EXEC FORTCLG,REGION.GO=200K *
    //FORT.SYSIN DD *
    MAIN PROGRAM

    //GO.SYSIN DD *
    DATA DECK
    /*

RRR000120
RRR000010
RRR000030
RRR000040
RRR000050
RRR000060
RRR000070
RRR000080
RPR000090
RRR000100
RRR000110
RRR000130
RRR000140
RRR000150
RRR000160
RRR000170
RRR000180
RRR000190
RRR000200
RRR000220
RRR000240
RRR000270
RRR000280
RRR000330
RPR000340
RRR000360
RRR000370
RRR000440
RRR000410
RRR000420
RPR000430
RRR000440
RRR000450
RRR000460
RRR000470
RRR000480

88

DATA INPUT DECK

THE FIRST CARD OF THE DATA DECK IS THE TITLE CARD.
IT MAY CONTAIN ANY ALPHA NUMERIC TITLE IN COLUMNS 1 - 80.

INTEGERS IN THE SECOND CARD IS THE PROBLEM CARD. IT CONTAINS
THE FIRST 13 COLUMNS IN THE FOLLOWING MANNER:
```
COL 1-4:  NUMBER OF OBSERVATIONS
COL 5-6:  NUMBER OF VARIABLES PER OBSERVATION
COL 7-8:  NUMBER OF CLUSTERS DESIRED
COL 9:    CRITERION TO BE USED IN THE EVALUATION:
          1 = TRACE W
          2 = DETERMINANT W
          3 = LARGEST ROOT OF W(INVERSE)*B
          4 = TRACE W(INVERSE)*B

CCL 10:   STANDARDIZATION PARAMETER:
          C = THE PROGRAM WILL NOT STANDARDIZE DATA
          1 = THE PROGRAM WILL STANDARDIZE DATA

CCL 11:   DISTANCE PARAMETER:
          0 = EUCLIDIAN DISTANCE
          1 = SCALED EUCLIDIAN DISTANCE
          2 = MAHALANOBIS DISTANCE

CCL 12:   DATA PARAMETER:
          C = DATA IS TO BE READ IN
          1 = NO DATA READ IN; PREVIOUS DATA IS
              TO BE REANALYZED USING CURRENT PROBLEM

CCL 13:   TIMING PARAMETER. OBSERVATIONS CONSIDERED IN
          0 = NO OBSERVATIONS USED INDIVIDUAL SWITCHES
          1 - 8 = AN INCREASING NUMBER OF OBSERVATIONS USED
          9 = ALL OBSERVATIONS USED
```

NOTE: THE TIMING PARAMETER DETERMINED HOW EXTENSIVE THE
DOUBLE CHECK OF THE CLUSTER SOLUTION IS TO BE.
NORMALLY "9" IS USED UNLESS THE DATA MATRIX IS LARGE

THIS CARD MUST BE IN (I4,2I2,5I1) FORMAT

THE THIRD CARD IS THE VARIABLE FORMAT CARD. THIS CARD
MUST BE PRESENT UNLESS THE DATA PARAMETER IN THE PROBLEM
CARD IS "1". THE FORMAT STATEMENT IS THAT TYPE WHICH THE DATA
CARDS ARE IN. FOR EXAMPLE IF THE DATA CARD HAS DATA:
```
ABCD 1.2     1.55
```
STARTING IN COL 1, THE VARIABLE FORMAT CARD WOULD
CONTAIN THE FOLLOWING STARTING IN COL 1:
```
(1A4,4X,F3.1,4X,F4.2)
```

THE NEXT CARDS ARE THE DATA CARDS. THEY MAY CONTAIN
DATA IN COLUMNS 1 - 72.

```
C   FOLLOWED FOR REANALYSIS OF THE SAME DATA, WITH A '1' IN COLUMN 11      RRRO0570
C   OF THE NEW PROBLEM CARD.                                              RRRO0580
C   SET INCLUDING NEW TITLE CARD, PROBLEM CARD, SIMPLY PLACE THE NEW DATA RRRO0590
C   AFTER THE PREVIOUS SET                                                RRRO1010
C                                                                         RRRO1020
C   FOR THE PROGRAM TO EXIT NORMALLY, TWO BLANK CARDS MUST FOLLOW         RRRO1030
C   THE LAST DATA BATCH.                                                  RRRO1040
C                                                                         RRRO1050
C     COMMON NOBS,NVARS,NGPS,ICRIT,NOSTAN,IDIST,IFINE,KTIME,IDENT(600),   RRRO1060
C    1DATA(600,20),T(20,20),B(20,20),WFCT(20,20),SVCEN(20,20),BT(20,20),  RRRO1070
C    2ICATA(600),NISV(20),VMEAN(20),XVEC(20,20),YVEC(20,20),BT(20,20),    RRRO1080
C    3NISVT(20),SVCENT(20,20),IDATAT(600,20),VEC(20,20),BIG(20),          RRRO1090
C                                                                         RRRO1100
C   DESCRIPTION OF COMMON AREA:                                           RRRO1110
C                                                                         RRRO1120
C   NOBS   = NUMBER OF OBSERVATIONS                                       RRRO1130
C   NVARS  = NUMBER OF VARIABLES                                          RRRO1140
C   NGPS   = NUMBER OF CLUSTERS                                           RRRO1150
C   ICRIT  = CRITERION TO BE OPTIMIZED (SEE ABOVE)                        RRRO1160
C   NOSTAN = STANDARDIZATION PARAMETER (SEE ABOVE)                        RRRO1170
C   IDIST  = DISTANCE PARAMETER (SEE ABOVE)                               RRRO1180
C   IFINE  = ESCAPE PARAMETER: IF IFINE IS SET EQUAL TO '1',              RRRO1190
C            SOMETHING IS WRONG AND THE APPROPRIATE ERROR                 RRRO1200
C            MESSAGE IS PRINTED OUT; THE PROGRAM GOES ON TO THE           RRRO1210
C            NEXT PROBLEM                                                 RRRO1220
C   KTIME  = TIMING PARAMETER (SEE ABOVE)                                 RRRO1230
C   IDENT  = OBSERVATION IDENTIFICATION (SEE ABOVE)                       RRRO1240
C            EACH DATA CARD                                               RRRO1250
C   DATA   = AREA IN WHICH THE DATA MATRIX IS STORED                      RRRO1260
C   T      = CROSS-PRODUCTS MATRIX                                        RRRO1270
C   B      = BETWEEN-CLUSTERS MATRIX                                      RRRO1280
C   W      = WITHIN-CLUSTERS MATRIX                                       RRRO1290
C   WFCT   = CHOLESKY FACTOR OF THE WITHIN-CLUSTERS MATRIX                RRRO1300
C   SVCEN  = CLUSTER CENTERS (MEANS)                                      RRRO1310
C   IDATA  = CLUSTER IDENTIFICATION FOR EACH OBSERVATION                  RRRO1320
C   NISV   = CLUSTER SIZES (NUMBER OF OBSERVATIONS IN EACH CLUSTER)       RRRO1330
C   VMEAN  = VARIABLE MEANS                                               RRRO1340
C   SD     = VARIABLE STANDARD DEVIATIONS                                 RRRO1350
C   XVEC   = TEMPORARY STORAGE                                            RRRO1360
C   YVEC   = TEMPORARY STORAGE                                            RRRO1370
C                                                                         RRRO1380
C   NAN=24                                                                RRRO1390
C   BT     = TEMPORARY STORAGE FOR BETWEEN-CLUSTERS MATRIX                RRRO1400
C   NISVT, SVCENT, IDATAT: TEMPORARY STORAGE SERVING THE                  RRRO1410
C            SAME FUNCTIONS AS NISV, SVCEN, AND IDATA                     RRRO1420
C                                                                         RRRO1430
```

```
C     VEC = EIGENVECTORS                                              RRRO1440
C     EIG = EIGENVALUES                                               RRRO1460
C                                                                     RRRO1470
  100 CALL PRELIM (CRIT)                                              RRRO1480
      IFINE IS ESCAPE PARAMETER                                       RRRO1490
      IF (IFINE.EC.1) GC TO 400                                       RRRO1500
      CALL RANDST (CRIT)                                              RRRO1510
      IF (IFINE.EC.2) GO TO 500                                       RRRO1520
  300 CALL KMEANS (CRIT)                                              RRRO1530
      IF (IFINE.EC.1) GC TO 400                                       RRRO1540
      CALL ISWTCH (CRIT)                                              RRRO1550
      GO TO 100                                                       RRRO1560
C     IFINE GOT SET TO 1:   IF DATA HAS BEEN READ IN, RESET DATA AND  RRRO1570
C                   NEXT PROBLEM                                      RRRO1580
      DO 410 I=1,NOBS                                                 RRRO1590
  400 DO 410 J=1,NVARS                                                RRRO1600
      IF (NOSTAN.EQ.1) DATA(I,J) = DATA(I,J) * SC(J)                  RRRO1610
  410 DATA(I,J) = DATA(I,J) + VMEAN(J)                                RRRO1630
      GO TO 100                                                       RRRO1640
  500 WRITE (6,900)                                                   RRRO1650
  900 FORMAT (' 0 END OF ANALYSES.')                                  RRRO1660
      STOP                                                            RRRO1670
      END                                                             RRRO1680
      SUBROUTINE PRELIM (CRIT)                                        RRRO1690
C                                                                     RRRO1700
C     THIS SUBROUTINE MAKES THE PRELIMINARY CALCULATIONS.             RRRO1710
C     IT INPUTS THE DATA, CALCULATES THE MEANS AND VARIANCES FOR EACH RRRO1720
C     VARIABLE, STANDARDIZES (CONVERTS TO Z-SCORES) EACH VARIABLE IF  RRRO1740
C     REQUESTED, AND CALCULATES THE CROSS-PRODUCTS MATRIX.            RRRO1750
C                                                                     RRRO1760
      COMMON NOBS,NVARS,NGPS,ICRIT,NOSTAN,IDIST,IFINE,KTIME,ICENT(600),RRRO1770
     1DATA(600,20),T(20,20),B(20,20),W(20),WFCT(20,20),SVCEN(20,20),   RRRO1780
     2NISVT(20),SVCENT(20,20),VMEAN(20),SC(20),IDATAT(600),VVEC(20),ET(20,20), RRRO1790
     3NISVT(20),IFMT(20),VAR(20),YVEC(20,20),VEC(20,20),EIG(20)        RRRO1800
      DIMENSION TITLE(20),IFMT(20),VAR(20)                            RRRO1810
C                                                                     RRRO1820
C     INPUT SECTION: READ IN TITLE CARD, PROBLEM CARD, OPTIONAL CARDS, RRRO1830
C     FORMAT CARD, AND DATA CARDS                                     RRRO1840
C     WRITE OUT SOLUTION SPECIFICATIONS                              RRRO1850
C                                                                     RRRO1860
      IFINE = 0                                                       RRRO1870
      READ (5,900) (TITLE(I),I=1,20)                                  RRRO1880
      WRITE (6,902) (TITLE(I),I=1,20)                                 RRRO1890
      READ (5,901) NOBS,NVARS,NGPS,ICRIT,NOSTAN,ICIST,NODATA,KTIME    RRRO1900
      IF (NOBS.EQ.0) GO TO 820                                        RRRO1910
```

91

```
      IF (NOBS.GT.600) GC TC 815                          RRRO1520
      IF (NVARS.GT.20) GC TC 815                          RPRO1530
      IF (ICRIT.GT.4) GC TO 815                           RRRO1540
      IF (NOSTAN.GT.1) GC TC 815                          RRRO1560
      IF (IDIST.GT.2) GC TC 815                           RRRO1570
      IF (NVARS.EC.1) ICRIT=1                             RRRO1580
    4 IF (NODATA.EQ.1) ICIST=0                            RRRO1590
    5 READ (5,902) (IFMT(I),I=1,18)                       RRRO2000
      WRITE (10,904) NOBS,NVARS,ICRIT                     RRRO2020
   10 WRITE (6,920)                                       RRRO2030
   20 WRITE (6,921)                                       RRRO2040
   30 WRITE (6,922)                                       RRRO2050
   40 WRITE (6,923)                                       RRRO2060
   50 CCNTINUE                                            RRRO2070
      IF (NCSTAN) 60,60,70                                RRRO2080
   60 WRITE (6,924)                                       RRRO2090
   70 WRITE (6,925)                                       RRRO2100
   80 CCNTINUE                                            RRRO2110
      IF (IDIST) 85,85,90                                 RRRO2120
   85 WRITE (6,926)                                       RRRO2130
   90 GC TO 95                                            RRRO2140
      IF (IDIST.EC.2) GC TC 92                            RRRO2150
   92 WRITE (6,936)                                       RRRO2160
   95 GC TO 95                                            RRRO2170
   96 CCNTINUE                                            RRRO2180
      IF (NODATA.EQ.1) GC TO 101                          RRRO2190
      WRITE (6,906) (IFMT(I),I=1,18)                      RRRO2200
      WRITE (6,980)                                       RRRO2210
  100 I=1,NOBS                                            RRRO2220
      REAC (5,IFMT) IDENT(I),(DATA(I,J),J=1,NVARS)        RRRO2230
      WRITE (6,950) IDENT(I),(DATA(I,J),J=1,NVARS)        RRRO2240
  100 CCNTINUE                                            RRRO2250
      NCCATA=1                                            RRRO2270
  101 GC TO 102                                           RRRO2290
  102 WRITE (6,943) KTIME                                 RRRO2300
C                                                         RRRO2310
C     CALCULATE VARIABLE MEANS ANC VARIANCES             RRRO2320
C     SLBTRACT CLT OVERALL MEAN                          RRRO2330
C                                                         RRRO2340
C                                                         RRRO2350
C                                                         RRRO2360
C                                                         RRRO2370
```

```
C
      DO 105 J=1,NVARS
      VMEAN (J) = 0.0
105   VAR (J) = 0.0
      DO 110 I = 1,NOBS
      DO 110 J = 1,NVARS
110   VMEAN (J) = VMEAN (J) + DATA(I,J)/NOBS
      DO 120 I=1,NOBS
      DO 120 J=1,NVARS
120   DATA(I,J) = DATA(I,J)-VMEAN(J)
      DO 125 J=1,NVARS
      VAR(J) = VAR(J) + (DATA(I,J)**2)/(NOBS-1)
      IF (VAR(J).LE.0.000001) MFLAG=1
125   SD(J) = SQRT (VAR(J))
      IF (MFLAG.EQ.1) GO TO 825
C
C     CALCULATE T = X'X, THE CROSS-PRODUCTS MATRIX
C
131   DO 135 K=1,NVARS
135   DO 135 J=K,NVARS
135   T(K,J) = 0.0
      DO 140 K=1,NVARS
      DO 140 J=K,NVARS
      DO 140 I=1,NOBS
140   T(K,J) = T(K,J) + DATA(I,K)*DATA(I,J)
C
C     OUTPUT SECTION: WRITE OUT MEANS, STANDARD DEVIATIONS, AND CROSS-
C                     PRODUCTS MATRIX
C
      WRITE (6,906) (VMEAN(J),J=1,NVARS)
      WRITE (6,907) (SD(J),J=1,NVARS)
      WRITE (6,911) (T(J,I),J=1,I)
      DO 810 I=1,NVARS
810   WRITE (6,955) (T(J,I),J=1,I)
      CONTINUE
      IF (NOSTAN.EQ.0) GO TO 180
C
C     STANDARDIZE IF REQUESTED
C
      DO 160 J=1,NVARS
      DO 150 I=1,NOBS
150   DATA(I,J) = DATA(I,J) / SD(J)
      DO 160 K=J,NVARS
160   T(J,K) = (1.0/SD(J))*T(J,K)*(1.0/SD(K))
      WRITE (6,702)
      DO 700 I=1,NOBS
      WRITE (6,701) (DATA(I,L),L=1,NVARS)
```

93

```
700 CONTINUE                                                    RRRO2830
701 FORMAT('0',10(F11.3,1X))                                    RRRO2840
702 FORMAT('1','STANDARDIZED DATA MATRIX IS')                   RRRO2850
    WRITE(6,912)                                                RRRO2860
    DO 270 I=1,NVARS                                            RRRO2870
270 WRITE(6,701)(T(J,I),J=1,I)                                  RRRO2880
17C CONTINUE                                                    RRRO2890
91C FORMAT('0 THE STANDARDIZED CROSS-PROD MATRIX IS ')          RRRO2900
81E RETURN                                                      RRRO2910
    IFINE=1                                                     RRRO2920
    IF(NODATA.EQ.0) IFINE=2                                     RRRO2930
    WRITE(6,935)                                                RRRO2940
    RETURN                                                      RRRO2950
82C IFINE = 2                                                   RRRO2960
    RETURN                                                      RRRO2970
825 WRITE(6,945) J                                              RRRO2980
    IFINE=1                                                     RRRO2990
    RETURN                                                      RRRO3010
                                                                RRRO2020
C     FORMAT STATEMENTS                                         RRRO3050
900 FORMAT (20A4)                                               RRRO3060
901 FORMAT (I4,2I2,5I1)                                         RRRO3070
902 FORMAT (18A4)                                               RRRO3080
9C4 FORMAT(' ',20A4)                                            RRRO3090
    1VARIABLES IS ',I2,/,' THE NUMBER OF OBSERVATIONS IS ',I4,/,' THE NUMBER OF   RRRO3100
905 FORMAT('0',/,' THE INPUT FORMAT IS ',18A4)                  RRRO3110
    GROUPS (INITIAL) IS ',I2)                                   RRRO3120
9C6 FORMAT('0',5(F11.3,1X),' MEANS ARE')                        RRRO3130
9C7 FORMAT('1',' THE VARIABLE STANDARD DEVIATIONS ARE')         RRRO3140
9C8 FORMAT('0THE VARIABLE STANDARD DEVIATIONS ARE')             RRRO3150
9C9 FORMAT('0OTHE CROSS-PRODUCTS MATRIX IS')                    RRRO3160
91C FORMAT('0OOTHE CRITERION IS TO MINIMIZE TRACE W')           RRRO3170
911 FORMAT('0OOTHE CRITERION IS TO MINIMIZE DETERMINANT W')     RRRC3190
912 FORMAT('0OOTHE CRITERION IS TO MAXIMIZE THE LARGEST ROOT OF W-1B')
913 FORMAT('0OOTHE CRITERION IS TO MAXIMIZE THE TRACE OF W-1B')
914 FORMAT('0OOTHE VARIABLES WILL NOT BE STANDARDIZED')
915 FORMAT('0OOTHE VARIABLES WILL BE STANDARDIZED')
916 FORMAT('0OOELCLIDIAN DISTANCE WILL BE USED')
917 FORMAT('0OOMAHALANOBIS DISTANCE WILL BE USED')
918 FORMAT('0OOERROR IN PROBLEM CARD, SKIPPING PROBLEM')
919 FORMAT('0OOWEIGHTED EUCLIDIAN DISTANCE WILL BE USED')
94C FORMAT('1','PREVIOUS DATA BEING USED')
945 FORMAT('0',/,' THE TIMING PARAMETER FOR INDIVIDUAL SWITCHES IS ',I1)
    1' WARNING: VARIABLE ',I2,' HAS ZERO VARIANCE.')
94E FORMAT('0',19.5X,10F8.3)
95C FORMAT('1',' INITIAL DATA MATRIX ',//)
    FORMAT('0',10(F10.2,1X))
```

94

```
      ENC
      SUBROUTINE RANDST (CRIT)                                          RRRC0320C
C                                                                       RRRC0321C
C     THIS SUBROUTINE DETERMINES THE INITIAL CLUSTER CONFIGURATICN.     RRRC0322C
C                                                                       RRRC0323C
      COMMON NOES,NVARS,NGPS,ICRIT,NOSTAN,IDIST,IFINE,KTIME,ICENT(6CC), RRRC0324C
     1CATA(6CO,20),T(20,20),B(20,20),W(20,20),WFCT(20,20),SVCEN(20,20), RRRC0325C
     2TCATA(6CO),AISV(20),VMEAN(20),SE(2C),XVEC(2C),YVEC(2C),ET(2C,2C),  RRRC0326C
     3 AISVT(20),SVCENT(20,20),IDATAT(6CC),VEC(20,2C),EIC(2C,2C)         RRRC0327C
      DIMENSICN ISTRT(2C)                                               RRRC0328C
C                                                                       RRRC0329C
C     ONE PASS K-MEANS TC DETERMINE INITIAL CLUSTER SOLUTION            RRRC0331C
C                                                                       RRRC0332C
C     RANDCM NUMBER STARTER                                             RRRC0334C
      IX=15739                                                          RRRC0335C
 200  NGFST=NGPS                                                        RRRC0336C
C 260: GETTING THREE INITIAL CONFIGURATICNS; THE CNE WITH THE           RRRC0370C
C     BEST CRITERICN VALUE WILL BE USED                                 RRRC0380C
 260  IBSRT = 1                                                         RRRC0390C
 205  M = 1,NGPS                                                        RRRC034CC
      ISVT(M)=1                                                         RRRC0410C
C     FIND A RANDOM OBSERVATION TO SERVE AS INITIAL CLUSTER CENTER      RRRC0420C
C     FOR EACH CLUSTER                                                  RRRC0430C
 201  RAND = URAN((IX)                                                  RRRC0440C
      ISTRT(M) = RAND * NOBS                                            RRRC0445C
      IF (ISTRT(M).EQ.0) ISTRT(M)=1                                     RRRC0460C
      IF (M.EC.1) GO TO 203                                             RRRC0470C
      MM1 = M-1                                                         RRRC0480C
      CO 204 MM=1,MM1                                                   RRRC0490C
      IF(ISTRT(M).EQ.ISTRT(MM)) GG TC 201                               RRRC0500C
 204  CONTINUE                                                          RRRC0510C
 203  HC = ISTRT(M)                                                     RRRC0520C
 205  CO 205 J=1,NVARS                                                  RRRC0530C
      SVCENT(M,J) = DATA(I,J)                                           RRRC0540C
      ITEMP = ICIST                                                     RRRC0550C
      ICIST = 0                                                         RRRC0560C
C     FOR EACH CBSERVATION                                              RRRC0570C
C     FIND EUCLIDIAN DISTANCE TO EACH CLUSTER CENTER                    RRRC0580C
 225  CO 215 M=1,NGPS                                                   RRRC0590C
 215  CO 210 J=1,NVARS                                                  RRRC0600C
      XVEC(J)=DATA(I,J)                                                 RRRC0610C
 210  YVEC(J) = SVCENT(M,J)                                             RRRC0620C
      CALL EISTCE (XVEC,YVEC,DISTA)                                     RRRC0630C
      IF (IFINE.EC.1) RETURN                                           RRRC0640C
      IF (M.EC.1) GO TO 212                                             RRRC0650C
      IF(DISTA.GE.SMDIST) GO TO 215                                     RRRC0660C
 212  SMDIST = DISTA                                                    RRRC0670C
```

95

```
215   ICATAT(I) = M
C     ASSIGN OBSERVATION TO CLOSEST CLUSTER AND UPDATE THAT CLUSTER
C     CENTER
      K=IDATAT(I)
      NISVT(K) = NISVT(K) + 1
225   CONTINUE J = 1,NVARS
      SVCENT(K,J) = ((NISVT(K)-1)*SVCENT(K,J)+DATA(I,J))/NISVT(K)
      IHIST = ITEMP
C     END OF DO FOR EACH OBSERVATION
C     RECALCULATE CLUSTER CENTERS TO ELIMINATE INITIAL RANDOM OBSERVA-
      TIONS
      DO 230 M = 1,NGPS
      NISVT(M) = NISVT(M)-1
230   CONTINUE J = 1,NVARS
      SVCENT(M,J) = 0.
      DO 235 I = 1,NOBS
      M = IDATAT(I)
      DO 235 J = 1,NVARS
235   SVCENT(M,J) = SVCENT(M,J) + DATA(I,J)/NISVT(M)

C     CALCULATE THE B AND W MATRICES
C     CALCULATE THE CRITERION VALUE

      ICIR=1
      IF (ICRIT.EQ.1) ICIR=1
      CALL WCALC (SVCENT,NISVT,NGPST,IDIR)
      IF (IFINE.EQ.1) RETURN
      CALL CRITON (CRIT)
      IF (IFINE.EQ.1) RETURN
C     WHICH INITIAL CONFIGURATION
C     FIRST INITIAL CONFIGURATION: BCRIT IS THE BEST CRITERION
      IF (IBSRT.GT.1) GO TO 250
C     FIRST INITIAL CONFIGURATION
245   BCRIT = CRIT
C     SAVE CLUSTER SIZES (NISV), CLUSTER CENTERS (SVCEN), CLUSTER
C     LISTS (LSTSV), AND OBSERVATION ID'S (IDATA)
      DO 247 M = 1,NGPS
      NISV(M) = NISVT(M)
247   DO 247 L = 1,NVARS
      SVCEN(M,L) = SVCENT(M,L)
      DO 249 I = 1,NOBS
249   IDATA(I) = IDATAT(I)
      GO TO 260
C     SECOND OR THIRD INITIAL CONFIGURATION
250   IF (CRIT.GE.BCRIT) GO TO 260
255   BCRIT = CRIT
      DO 257 M = 1,NGPS
      NISV(M) = NISVT(M)
```

```
 257  DC 257 L = 1,NVARS
      SVCEN(M,L) = SVCEN(M,L)
 259  DC 259 I=1,NCBS
 260  ICATA(I) = ICATAT(I)
      CCNTINUE
      CALL CCNE
      RETURN
      EAC
C     SLEROUTINE KMEANS (CRIT)
C
C     K-MEANS.  EACH OBSERVATION IS ASSIGNED TC THE CLOSEST CLUSTER
C     CENTER..
C
      CCMMCN NOBS,NVARS,NGPS,ICRIT,NOSTAN,IDIST,IFINE,KTIME,IDENT(600);
     1ICATA(600,20),T(20,20),B(20,20),SD(20),VMEAN(20,20),SVCEN(20,20);
     2ICATA(600),NISV(20),SVCENT(20,20),IDATAT(600),XVEC(20,20),BT(20,20);
     3NISVT(20),SVCENT(20,20),VEC(20,20),YVEC(20,20),EIG(20)
      INITIALIZING FLAGS
 200  ANC TEMPCRARY STORAGE FOR NISV, SVCEN
      NGPST = NGPS
      JFLAG = C
      ECRIT = CRIT
      DC 201 M=1,NGPS
      NISVT(M) = NISV(M)
 201  SVCENT(M,J) = SVCEN(M,J)
C     RECALCULATING WITHIN-CLUSTERS MATRIX
      ICIR=1
      IF (IDIST.EC.2) ICIR=2
      CALL WCALC (SVCENT,NISVT,NGPST,ICIR)
      IF (IFINE.EC.1) RETURN
C     INITIALIZING TEMPORARY STORAGE FCR LSTSV, IDATA
 210  DC 210 I = 1,NOBS
      ICATAT(I) = IDATA(I)
C     MAJOR DO LOOP: DO FOR EACH OBSERVATION
 400  DC 400 I = 1,NOBS
C
C     BEGIN K-MEANS SECTICN
C
C     CALCULATE VECTOR OF CISTANCES FRCM CBSERVATION TO EACH CLLSTER
 325  CENTER
      DC 335 M=1,NGPST
 330  DC 330 J=1,NVARS
      XVEC (J) = CATA (I,J)
 320  YVEC (J) = SVCENT (M,J)
      CALL DISTCE (XVEC,YVEC,DISTB)
      IF (IFINE.EC.1) RETURN
```

97

```
      IF (M.EC.1) GO TO 332
      IF (DISTB.GE.SMDIST) GO TO 335
      SMDIST = DISTB
      ICGP = M
332   CONTINUE
C
C     IS OBSERVATION ALREADY IN THE CLOSEST CLUSTER? IF YES, SKIP THIS
C     SECTION.
335   IF (ICGP.EQ.IDATAT(I)) GO TO 360
C
C     ICLD IS OLD CLUSTER ASSIGNMENT
346   ICLC = IDATAT(I)
C     INEW IS NEW CLUSTER ASSIGNMENT
      INEW = ICGP
C
      IFLAG = 1
C     RECALCULATE CLUSTER CENTERS
      DO 348 J=1,NVARS
      SVCENT(ICLD,J) = (NISVT(IOLD)*SVCENT(IOLD,J)-CATA(I,J))/
     1 (NISVT(IOLD)-1)
348   SVCENT(INEW,J) = (NISVT(INEW)*SVCENT(INEW,J)+CATA(I,J))/
     1 (NISVT(INEW)+1)
C     ADJUST NISV
      NISVT(IOLD) = NISVT(IOLD)-1
      NISVT(INEW) = NISVT(INEW)+1
C     ADJUST IDATA
      IDATAT(I) = IDGP
C     RECALCULATE WITHIN-CLUSTERS MATRIX FOR USE IN COMPUTING SCALED
C     EUCLICIAN AND MAHALANOBIS DISTANCE
      IF (IDIST.EC.0) GC TO 360
      CALL WCALC (SVCENT,NISVT,NGPST,ICIR)
      IF (IFINE.EC.1) RETURN
360   CONTINUE
C     DONE WITH MAJOR DO LOOP
400   CONTINUE
C
C     CALCULATE THE CRITERION
C
C     RECALCULATE SVCEN:    ACCURACY MEASURE
      DO 401 M=1,NGPST
      DO 401 J=1,NVARS
401   SVCENT(M,J)=0.
      DO 402 I=1,NOBS
      M = IDATAT(I)
      DO 402 J=1,NVARS
402   SVCENT(M,J) = SVCENT(M,J) + CATA(I,J)/NISVT(M)
C     RECALCULATE WITHIN-CLUSTERS MATRIX AND CRITERION VALUE BASED ON
C     NEW CLUSTER CENTER VALUES
      ICIR = 2
      IF (ICRIT.EC.1) ICIR = 1
      CALL WCALC (SVCENT,NISVT,NGPST,ICIR)
```

98

```
      IF (IFINE.EC.1) RETURN                                          RRRO5120
405   CALL CRITCH (CRIT)                                              RRRC513C
      IF (IFINE.EC.1) RETURN                                          RRRC5140
C     IF CRITERION BETTER THAN BEFORE? IF YES, THEN ANOTHER ITERATICN: RRRO5150
      IF (NO, FINISH KMEANS                                           RRRC5160
      IF (CRIT.GE.BCRIT) GO TO 535                                    RRRO517C
C                                                                     RRRO5180
C     ANOTHER ITERATION                                              RRRO519C
C                                                                     RRRC5200
515   PUT TEMPORARY VALUES INTO PERMANENT LOCATICNS                  RRRC521C
      NGPS=NGPST                                                     RRRO522C
      DO 520 M = 1,NGPS                                               RRRC523C
      NISV(M) = NISVT (M)                                            RRRC524C
520   DO 520 J = 1,NVARS                                             RRRC525C
      SVCEN (M,J) = SVCENT (M,J)                                     RRRC526C
530   DO 530 I = 1,NOBS                                              RRRO527C
      IDATA(I) = IDATAT(I)                                           RRRO528C
      GO TO 200                                                      RRRO529C
C                                                                     RRRO5300
C     FINISH                                                         RRRO531C
C                                                                     RRRO5132C
535   JFLAG= C  MEANS NO CHANGES HAVE BEEN MADE DURING THE LAST      RRRO513C
C     ITERATION: ITERATICNS CONVERGED                               RRRO5134C
      IF (JFLAG.EC.0) RETURN                                        RRRC5135C
C     JFLAG = 1 MEANS CHANGES HAVE BEEN MADE BUT THE CRITERICN VALUE RRRO5360
C     GOT WORSE: ITERATICNS NOT CONVERGING                          RRRO537C
540   WRITE (6,940) CRIT                                             RRRC538C
      WRITE (6,942) BCRIT                                            RRRO539C
      RECALCULATE WITHIN-CLUSTERS MATRIX AND RESET CRIT             RRRC54CC
      ICIR=2                                                        RRRO5410
      IF (ICRIT.EC.1) ICIR=1                                        RRRO542C
      CALL WCALC (SVCEN,NISV,NGPS,IDIR)                             RRRO543C
      IF (IFINE.EC.1) RETURN                                        RRRO5440
      CRIT=BCRIT                                                    RRRO5450
      RETURN                                                        RRRO546C
940   FORMAT ('0 ITERATIONS NOT CONVERGING')                       RRRO547C
942   FORMAT ('0 THE CRITERION VALUE IS ',E12.6)                   RRRC548C
943   FORMAT ('0 THE BEST CRITERION VALUE IS ',E12.6)              RRRC549C
      END                                                          RRRO550C
      SUBROUTINE ISWTCH (CRIT)                                     RRRO551C
C                                                                   RRRO5530
C     THIS SUBROUTINE CONSIDERS SWITCHING EACH OBSERVATION TO A    RRRO5540
C     DIFFERENT CLUSTER. THE SWITCH IS MADE IFF A BETTER CRITERION RRRO5550
C     VALUE RESULTS.                                               RRRO5560
C     THIS SUBROUTINE ALSO DEPENDS ON THE PARAMETER KTIME: IT DETERMINES RRRO5570
C     WHICH OBSERVATIONS ARE TO BE CONSIDERED. FOR COMPLETE EXPLANATICN RRRC5580
C     OF THE KTIME PARAMETER, SEE THE PROGRAM DESCRIPTICN.         RRRO5590
```

99

```fortran
      COMMON NOBS,NVARS,NGPS,ICRIT,NOSTAR,IDIST,IFINE,KTIME,IDENT(600),
     1DATA(600),T(20,20),B(20,20),W(20,20),SC(600),SVCEN(20,20),
     2ICATA(600),NISV(20),VMEAN(20),IDATAT(600),VEC(20,20),B1(20,20),
     3 NISVT(20),SVCENT(20),VCENT(20)
      DIMENSION JFLAG(20)
C     SET UP THE ICIR FLAG
      ICIR=2
      IF(ICRIT.EC.1) ICIR=1
C     KTIME = 0 MEANS SKIP THIS HEURISTIC
      IF(KTIME.EC.0) GC TO 750
C     SET UP TEMPORARY STORAGE AREAS FOR NGPS, NISV, AND SVCEN
      NGPST=NGPS
      DO 500 M=1,NGPS
      NISVT(M)=NISV(M)
      JFLAG(M)=0
  500 DO J = 1,NVARS
      SVCENT(M,J)=SVCEN(M,J)
C     JFLAG WILL BE SET = 1 IF ANY SWITCH IS MADE
  600 IFLAG=0
C     MAJOR DO LOOP:
      DO 700 M=1,NGPS
C     JFLAG(M) INDICATES WHETHER CLUSTER M HAS BEEN ALTERED
      IF(JFLAG(M).EQ.2) GO TO 700
C     KTIME = MEANS ALL OBSERVATIONS WILL BE CONSIDERED
      IF(KTIME.EC.9) GC TO 617
C     COMPUTE ALL DISTANCES FROM OBSERVATIONS TO CLUSTER CENTERS FOR THE
C     OBSERVATIONS IN CLUSTER M
      DO 605 J=1,NVARS
  605 YVEC(J) = SVCEN(M,J)
      BIGDIS = 0.
      DO 615 I=1,NOBS
      IF(IDATA(I).NE.M) GO TO 615
      DO 610 J=1,NVARS
  610 XVEC(J)=DATA(I,J)
      CALL DISTCE (XVEC,YVEC,DISTA)
C     LOCATE OBSERVATION WITH BIGGEST DISTANCE
      IF(BIGDIS.GE.DISTA) GO TO 615
      BIGDIS=CISTA
  615 CONTINUE
C     DIVIDE BIGDIS BY TIMING PARAMETER KTIME
      FSIVE = KTIME
      FSIVE = BIGDIS / TIME
C     FOR ALL OBSERVATIONS IN CLUSTER M
      DO 691 I=1,NOBS
      IF(IDATA(I).NE.M) GC TO 691
      KFLAG = 0
```

RRRO5600
RRRO5610
RRRO5620
RRRO5630
RRRO5640
RRRO5650
RRRO5660
RRRO5670
RRRO5680
RRRO5690
RRRO5700
RRRO5710
RRRO5720
RRRO5730
RRRO5740
RRRO5750
RRRO5760
RRRO5770
RRRO5780
RRRO5790
RRRO5800
RRRO5810
RRRO5820
RRRO5830
RRRO5840
RRRO5850
RRRO5860
RRRO5870
RRRO5880
RRRO5890
RRRO5900
RRRO5910
RRRO5920
RRRO5930
RRRO5940
RRRO5950
RRRO5960
RRRO5970
RRRO5980
RRRO5990
RRRO6000
RRRO6010
RRRO6020
RRRO6040
RRRO6050
RRRO6060
RRRO6070

100

```
C     KTIME = 9 MEANS CONSIDER ALL OBSERVATIONS                          RRRO6080C
C     IF (KTIME.EQ.9) GO TO 618                                          RRRO609C
C     CURRENT CLUSTER CENTER IS LT BIGDIS IF THE DISTANCE TO ITS         RRRO6100C
      DO 619 J=1,NVARS                                                   RRRO6110C
619   XVEC(J) = DATA(I,J)                                                RRRO612C
      CALL DISTCE (XVEC,YVEC,DISTA)                                      RRRO6130
      IF (DISTA.LT.PSME) GO TO 691                                       RRRO614C
618   CONTINUE                                                           RRRO6150
      DO 690 MNEW=1,NGPS                                                 RRRO6160
      IF (NISV(MOLD).EQ.1) GO TO 690                                     RRRO617C
      IF (MNEW.EQ.MOLD) GO TO 690                                        RRRO6180
      IF (KFLAG.EQ.1) GO TO 690                                          RRRO619C
C     COMPUTE NEW SVCEN BASED ON OBSERVATION BEING SWITCHED              RRRO6200
      DO 620 J=1,NVARS                                                   RRRO6210
      SVCENT(MNEW,J) = (NISV(MNEW)*SVCEN(MNEW,J)+DATA(I,J))/             RRRO6220
     1 (NISV(MNEW)+1)                                                    RRRO623C
      SVCENT(MOLD,J) = (NISV(MOLD)*SVCEN(MOLD,J)-DATA(I,J))/             RRRO624C
     1 (NISV(MOLD)-1)                                                    RRRO625C
620   CONTINUE                                                           RRRO626C
C     ADJUST NISV                                                        RRRO627C
      NISVT(MNEW) = NISV(MNEW)+1                                         RRRO6280
      NISVT(MOLD) = NISV(MOLD)-1                                         RRRO629C
C     COMPUTE WITHIN-CLUSTERS MATRIX AND NEW CRITERION VALUE             RRRO630C
      CALL WCALC (SVCENT,NISVT,NGPST,ICIR)                               RRRO6310
      IF (IFINE.EQ.1) RETURN                                            RRRO6320
      CALL CRITON (TCRIT)                                                RRRO6330
C     MAKE THE SWITCH IFF NEW CRIT IS BETTER                             RRRO6340
      IF (IFINE.EQ.1) RETURN                                            RRRO6350
      IF (TCRIT.GE.CRIT) GO TO 680                                       RRRO6360
                                                                         RRRO637C
C     MAKE THE SWITCH                                                    RRRO6380
                                                                         RRRO639C
C     RESET FLAGS                                                        RRRO640C
      JFLAG(M) = 1                                                       RRRO641C
      JFLAG(MNEW) = 1                                                    RRRO642C
      IFLAG=1                                                            RRRO6430
      KFLAG = 1                                                          RRRO644C
C     ADJUST CRIT, IDATA, NISV, SVCEN                                    RRRO645C
      CRIT=TCRIT                                                         RRRO646C
      IDATA(I) = MNEW                                                    RRRO647C
      NISV(MNEW) = NISVT(MNEW)                                           RRRO648C
      NISV(MOLD) = NISVT(MOLD)                                           RRRO649C
      DO 630 J=1,NVARS                                                   RRRO650C
630   SVCEN(MNEW,J) = SVCENT(MNEW,J)                                     RRRO651C
      SVCEN(MOLD,J) = SVCENT(MOLD,J)                                     RRRO6520
630   GO TO 690                                                          RRRO653C
C     SWITCH WAS NOT MADE; RESET NISVT, SVCENT, TC VALUES PRESENT        RRRO654C
                                                                         RRRO655C
```

```
C     BEFORE THE SWITCH WAS CONSIDERED
      NISVT(MNEW) = NISV(MNEW)                                    RRRO06560
      NISVT(MOLD) = NISV(MOLD)                                    RRRO06570
      DO 685 J=1,NVARS                                            RRRO06580
      SVCENT(MNEW,J) = SVCEN(MNEW,J)                              RRRO06590
      SVCENT(MOLD,J) = SVCEN(MOLD,J)                              RRRO06600
  685 CONTINUE                                                    RRRO06610
  690 CONTINUE                                                    RRRO06620
C                                                                 RRRO06630
C     FINISH WITH CLUSTER M:  IF NO SWITCHES HAVE BEEN MADE SET   RRRO06640
C     JFLAG(M) = 2 AND GO TO NEXT CLUSTER; IF SWITCHES HAVE BEEN MADE, RRRO06650
      IF (JFLAG(M).EQ.0) GO TO 699                                RRRO06660
      JFLAG(M) = C                                                RRRO06670
      GO TO 700                                                   RRRO06680
  699 JFLAG(M) = 2                                                RRRO06690
  700 CONTINUE                                                    RRRO06700
C     DONE WITH ALL CLUSTERS:  IFLAG = 1 MEANS SOME SWITCHES HAVE BEEN RRRO06710
C     MADE; GO BACK AND ITERATE                                   RRRO06720
      IF (IFLAG.EC.1) GO TO 600                                   RRRO06730
C                                                                 RRRO06740
C     ALL DONE; ACCURATELY CALCULATE MEANS AND CRITERIA AND OUTPUT RRRO06750
C     THE RESULTS                                                 RRRO06760
  750 WRITE (6,900)                                               RRRO06770
  900 FORMAT (*1*,'THE FINAL CLUSTER SOLUTION IS ')               RRRO06780
C     RECALCULATE CLUSTER CENTERS, WITHIN-CLUSTERS MATRIX, AND CRITERION RRRO06790
      VALUE                                                       RRRO06810
      DO 715 M=1,NGPS                                             RRRO06820
      DO 715 J=1,NVARS                                            RRRO06830
  715 SVCEN(M,J) = 0.                                             RRRO06840
      DO 720 I=1,NOBS                                             RRRO06850
      M = IDATA(I)                                                RRRO06860
      DO 720 J=1,NVARS                                            RRRO06870
  720 SVCEN(M,J) = (SVCEN(M,J) +DATA(I,J))/NISV(M)                RRRO06880
      CALL WCALC (SVCEN,NISV,NGPS,IDIR)                           RRRO06890
      IF (IFINE.EC.1) RETURN                                      RRRO06900
      CALL CRITION (CRIT)                                         RRRO06910
      IF (IFINE.EC.1) RETURN                                      RRRO06920
C     CALL THE OUTPUT ROUTINE                                     RRRO06930
      CALL OUTPUT (CRIT)                                          RRRO06940
      WRITE (6,901)                                               RRRO06950
  901 FORMAT (*0 END OF CLUSTER PROBLEM*)                         RRRO06960
      RETURN                                                      RRRO06970
      END                                                         RRRO06980
      SUBROUTINE OUTPUT (CRIT)                                    RRRO06990
C                                                                 RRRO07000
C     THIS SUBROUTINE PRINTS OUT THE CLUSTER SOLUTION             RRRO07010
C     FOR EACH CLUSTER - THE CLUSTER SIZE                         RRRO07020
C                                                                 RRRO07030
```

102

```
C          THE CLUSTER CENTROID
C          THE OBSERVATIONS BELONGING TO THE CLUSTER
C          THE WITHIN CELLS MATRIX
C          THE CRITERION VALUE
       COMMON NOBS,NVARS,NGPS,ICRIT,NOSTAN,IDIST,IFINE,KTIME,IDENT(600),
      1DATA(600,20),T(20,20),B(20,20),WC(20,20),EC(600),SVCEN(20,20),
      2NISVT(20),SVCENT(20),VMEAN(20),IDATAT(600),VEC(20,20),EIG(20)
C      UNSTANDARDIZE IF NECESSARY
       IF (NOSTAN.EQ.0) GO TO 40
       DO 20 M=1,NGPS
       DO 20 J=1,NVARS
   20  SVCEN(M,J) = SVCEN(M,J)*SD(J)
       DO 30 J=1,NVARS
       DO 30 K=J,NVARS
   30  T(J,K) = SD(J)*T(J,K)*SD(K)
       ICIR = 2
       IF (ICRIT.EQ.1) ICIR=1
       CALL WCALC (SVCEN,NISV,NGPS,IDIR)
       IF (IFINE.EQ.1) RETURN
       CALL CRITON (CRIT)
       IF (IFINE.EQ.1) RETURN
C      THE OVERALL MEAN BACK INTO EACH OBSERVATION: UNSTANDARDIZE
C      THE DATA IF NECESSARY
       DO 80 J=1,NVARS
       DO 50 I=1,NOBS
       IF (NOSTAR.EQ.1) DATA(I,J) = DATA(I,J)*SD(J)
   50  DATA(I,J) = DATA(I,J) + VMEAN(J)
       DO 60 M=1,NGPS
   60  SVCEN(M,J) = SVCEN(M,J) + VMEAN(J)
   80  CONTINUE
C      WRITE OUT NISV, SVCEN
       WRITE (6,90C) M
       WRITE (6,901) NISV(M)
       WRITE (6,902) (SVCEN(M,J),J=1,NVARS)
   55  WRITE (7,915) (SVCEN(M,J),J=1,NVARS)
  915  FORMAT (8F8.4)/8F8.4)
       WRITE (6,903) NISV(M)
       ISTNC = NISV(M)
       DO 100 I=1,NOBS
       IF (IDATA(I).NE.M) GO TO 100
       K=K+1
       IDATAT(K) = IDENT(I)
  100  CONTINUE
C      SORT THE OBSERVATION IDENTIFICATIONS
```

                                                    RRRC7C4C
                                                    RRRC07050
                                                    RRRC07C6C
                                                    RRRC07C7C
                                                    RRRO07080
                                                    RRRO07090
                                                    RRRC07100
                                                    RRRO07110
                                                    RRRO07120
                                                    RRRO07130
                                                    RRRO07140
                                                    RRRO07150
                                                    RRRO07160
                                                    RRRO07170
                                                    RRRO07180
                                                    RRRO07190
                                                    RRRC072CC
                                                    RRRO07210
                                                    RRRO07220
                                                    RRRO07230
                                                    RRRO07240
                                                    RRRO07250
                                                    RRRO07260
                                                    RRRO07270
                                                    RRRO07280
                                                    RRRO07290
                                                    RRRO07300
                                                    RRRO07310
                                                    RRRO07320
                                                    RRRO07330
                                                    RRRO07340
                                                    RRRO07350
                                                    RRRO07360
                                                    RRRO07370
                                                    RRRO07380
                                                    RRRO07390
                                                    RRRO074CC

                                                    RRRO07410
                                                    RRRC07420
                                                    RRRO07430
                                                    RRRO07440
                                                    RRRC07450
                                                    RRRO07460
                                                    RRRO07480
                                                    RRRO07490

```
      L=0
101   L=L+1
      IF(L.EQ.NISV(M)) GO TO 110
102   LL=LL+1
      IF((ICATAT(L).LE.ICATAT(LL)) GO TO 101
      LTEMP = IDATAT(L)
      ICATAT(L) = IDATAT(LL)
      ICATAT(LL) = LTEMF
      IF(L.EC.1) GO TO 102
      L=L-1
      GO TO 102
C
110   CCNTINUE
C     WRITE OUT THE IDENTIFICATIONS
      WRITE(6,904)(IDATAT(K),K=1,LSTNC)
      WRITE(7,920)(IDATAT(K),K=1,LSTNC)
920   FCRMAT(16I5)
115   CCNTINUE
C     WRITE OUT THE WITHIN-CLUSTERS MATRIX
      WRITE(6,905)
      DO 120 J=1,NVARS
      WRITE(6,906) (W(K,J),K=1,J)
120   CCNTINUE
C     ICRIT IS THE CRITERION CHOSEN BY THE USER
      GOTO (15C,160,17C),ICRIT
C     WRITE OUT TRACE W
150   WRITE(6,907) CRIT
      RETURN
C     WRITE OUT DET W
160   ANC LOG (DET T / DET W)
      CALL UPRFCT (NVARS,T,M)
      DET =1.
      DO 165 J=1,NVARS
165   DET = DET*T(J,J)
      WRITE(6,912) CRIT
      CRIT = ALOG 10 (((ET**2)/(CRIT**2))
      WRITE(6,90C) CRIT
      RETURN
C     FCF LARGEST FOOT ANC HOTELLING'S TRACE CRITERIA, FINC EIGENVALLES
      DO 175 J=1,NVARS
170   DO 175 K=J,NVARS
      BT(J,K) = B(J,K)
175   BT(K,J) = BT(J,K)
      CFIT = 1.0 / CRIT
      CALL UTISUI (NVARS,BT,WFCT)
      CALL EIGN (NVARS,ET,EIG,VEC,IND)
      IF(IND.NE.C) GO TO 180
C     WRITE OUT EIGENVALLES
```

```
RRRO7500
RRRO7510
RRRO7520
RRRO7530
RRRO7540
RRRO7550
RRRO7560
RRRO7570
RRRO7580
RRRO7590
RRRO7600
RRRO7620
RRRO7630
RRRO7640
RRRO7650
RRRO7660
RRRO7670
RRRO7680
RRRO7690
RRRO7700
RRRO7720
RRRO7730
RRRO7740
RRRO7750
RRRO7760
RRRO7770
RRRO7780
RRRO7790
RRRO7800
RRRO7820
RRRO7830
RRRO7840
RRRO7850
RRRO7870
RRRO7880
RRRO7890
RRRO7900
RRRO7910
RRRO7920
RRRO7930
RRRO7940
RRRO7950
```

104

```
      WRITE (6,910) (EIG(J),J=1,NVARS)
      IF(ICRIT.EQ.4)GO TO 178                                           RRR07560
C                                                                       RRR07570
C     WRITE OUT LARGEST FCCT CRIT                                       RRR07580
      RETURN                                                            RRR07590
  178 WRITE (6,914) CRIT                                                RRR08000
      RETURN                                                            RRR08010
C                                                                       RRR08020
C     ESCAPE ROUTE                                                      RRR08030
  180 IFINE = 1                                                         RRR08040
      WRITE (6,911) IND                                                 RRR08050
      RETURN                                                            RRR08060
C                                                                       RRR08070
C     FORMAT STATEMENTS                                                 RRR08080
C                                                                       RRR08090
  900 FORMAT ('CCLUSTER ',I2)                                           RRR08100
  901 FORMAT ('  SIZE = ',I3)                                           RRR08110
  902 FORMAT ('  CENTER = ',I3)                                         RRR08120
  903 FORMAT (' THE OBSERVATIONS ARE ')
  904 FORMAT (8X,5I7,/,5X,5I7)
  905 FORMAT (' THE WITHIN GROUP MATRIX IS ')                           RRR0814C
  906 FORMAT (10(F10.3,1X))
  907 FORMAT (' THE TRACE OF THE WITHIN-CLUSTERS MATRIX IS ',E12.6)     RRR08170
  908 FORMAT (' CETERMINANT T IS ',E12.6)                               RRR0818C
  909 FORMAT (' THE EIGENVALUES OF W(INVERSE)*E ARE ',10F8.3)           RRR08190
  910 FORMAT (' LOG (DET)/ DET W(INVERSE)-INC IS ',I3)                  RRR08200
  911 FORMAT (' ILL-CCNDITIONED W MATRIX --')                           RRR08210
  912 FORMAT (' THE DETERMINANT OF THE WITHIN-CLUSTERS MATRIX IS ',E12.6)RRR08230
  913 FORMAT (' THE LARGEST RCCT IS ',E12.6)                            RRR08240
  914 FORMAT (' THE TRACE OF W(INV)*B IS ',E12.6)                       RRR08250
      END                                                               RRR08260
      SUBROUTINE DISTCE (XVECS,YVECS,CIST)                              RRR08270
C                                                                       RRR08280
C     THIS SUBROUTINE CALCULATES EUCLICIAN, WEIGHTED EUCLICIAN, CR      RRR08290
C     MAHALANOBIS DISTANCE (DEPENDING CR IDIST).                        RRR08300
C                                                                       RRR08310
      COMMON NOBS,NVARS,NGPS,ICRIT,NCSTAN,IDIST,IFINE,KTIME,ICENT(600), RRR08320
     1DATA(600,20),T(20,20),B(20,20),WFCT(20,20),SVCEN(20,20),          RRR08330
     2ICATA(600),AISV(20),VMEAN(20),SE(20),IDATAT(600),XVEC(20,20),YVEC(20,20), RRR08340
     3AISVT(20),SVCENT(20),IDATAT(600),XVECS(1),YVECS(1)                RRR08350
      DIMENSION X(20),XVECS(1),YVECS(1)                                 RRR08360
      ICIR = IDIST+1                                                    RRR08370
      GO TO (100,200,300),IDIR                                          RRR08390
  900 WRITE (6,900)                                                     RRR08400
      FORMAT (' ERROR IN CISTANCE CODE - COLUMN 10 CF PRCELEM CARD')    RRR08410
      IFINE = 1                                                         RRR08420
      RETURN                                                            RRR08430
C
```

105

```
C     CALCULATE EUCLIDIAN DISTANCE                                RRRO0844C
100   DIST=0.                                                     RRRO0845C
      DO 150 J=1,NVARS                                            RRRO0846C
150   DIST = DIST + (XVECS(J)-YVECS(J))**2                        RRRO0847C
      RETURN                                                      RRRO0848C
C                                                                 RRRO0849C
C     CALCULATE EUCLIDIAN DISTANCE FOR STANDARDIZED VARIABLES     RRRO0850C
200   DIST=0.                                                     RRRO0851C
      DO 205 J=1,NVARS                                            RRRO0852C
205   DIST = DIST + (XVECS(J)-YVECS(J))*(1.0/W(J,J))*(XVECS(J)-YVECS(J)) RRRO0853C
      RETURN                                                      RRRO0854C
C                                                                 RRRO0855C
C     CALCULATE MAHALANOBIS DISTANCE AS THE ELEMENTS OF           RRRO0856C
C     L'(INV)(XVEC-YVEC) SQUARED WHERE L IS THE CHOLESKY FACTOR OF W RRRO0857C
300   DIST=0.                                                     RRRO0858C
      DO 305 J=1,NVARS                                            RRRO0859C
305   X(J) = XVECS(J) - YVECS(J)                                  RRRO0860C
      CALL UTIRT(NVARS,1,WFCT,X)                                  RRRO0861C
      DO 310 J=1,NVARS                                            RRRO0862C
310   DIST = DIST + X(J)**2                                       RRRO0863C
      RETURN                                                      RRRO0864C
      END                                                         RRRO0865C
      SUBROUTINE CRITON (CRIT)                                    RRRO0866C
C                                                                 RRRO0867C
C     THIS SUBROUTINE CALCULATES THE CRITERIA VALUES              RRRO0868C
      COMMON NOBS,NVARS,NGPS,ICRIT,NOSTAN,IDIST,IFINE,KTIME,IDENT(600), RRRO0869C
     1DATA(600,20),T(20,20),B(20,20),W(20,20),WFCT(20,20),SVCEN(20,20), RRRO0870C
     2IDATA(600),AISV(20),VMEAN(20),SC(20),XVEC(20),VEC(20,20),ET(20,20), RRRO0871C
     3AISVT(20),SVCENT(20,20),IDATAT(600),VEC(20,20),EIG(20)      RRRO0872C
      GO TO (100,200,300,400),ICRIT                               RRRO0873C
      WRITE (6,90)                                                RRRO0874C
90    FORMAT (6/1 ERROR IN CRITERION CODE, COLUMN 6 IN PROBLEM CARD.') RRRO0875C
      IFINE =1                                                    RRRO0876C
      RETURN                                                      RRRO0877C
C                                                                 RRRO0878C
C     CALCULATE TRACE W                                           RRRO0879C
100   CRIT = 0.                                                   RRRO0880C
      DO 105 J=1,NVARS                                            RRRO0881C
105   CRIT = CRIT + W(J,J)                                        RRRO0882C
      RETURN                                                      RRRO0883C
C                                                                 RRRO0884C
C     CALCULATE DET W AS THE PRODUCT OF DIAGONAL ELEMENTS OF THE  RRRO0890C
C     CHOLESKY FACTOR OF W                                        RRRO0891C
```

106

```
C
      DET = 1.
2CC   DO 205 J=1,NVARS
205   DET = DET*WFCT(J,J)
      DFIT = DET
      RETURN
C
C     CALCULATE LARGEST ROOT AS L(INV)*B*L'(INV) WHERE L IS THE
C     CHOLESKY FACTOR OF W
C     CRITERION IS RECIPROCAL OF LARGEST ROOT
300   DO 305 J=1,NVARS
      DO 305 K=J,NVARS
      BT(J,K) = B(J,K)
305   BT(K,J) = BT(J,K)
      CALL EIGN (NVARS,BT,WFCT)
      CALL UTISUI (NVARS,BT,EIG,VEC,IND)
      CFIT = 1.0 / EIG(1)
      RETURN
C
C     CALCULATE HOTELLING'S TRACE AS SUM OF DIAGONAL ELEMENTS OF
C     L(INV)*B*L'(INV)
C     CRITERION IS RECIPROCAL OF THIS SUM
4CC   DO 405 J=1,NVARS
      DO 405 K=J,NVARS
      BT(J,K) = B(J,K)
4C5   BT(K,J) = BT(J,K)
      CALL UTISUI (NVARS,BT,WFCT)
      CFIT = 0.
      DO 410 J=1,NVARS
41C   CFIT = CFIT+BT(J,J)
      CFIT=1.0/CFIT
      RETURN
      END
      SUBROUTINE WCALC (SV,NI,NGP,IGIR)
C
C     THIS SUBROUTINE CALCULATES THE WITHIN-CELLS MATRIX AND, IF
C     NECESSARY, THE CHOLESKY FACTOR OF THE WITHIN-CELLS MATRIX
C
      COMMON NOES,NVARS,NGPS,ICRIT,NOSTAN,IDIST,IFINE,KTIME,ICENT(6CC),
     1IDATA(6C0,20),T(20,20),B(20,20),WW(20,20),WFCT(20,20),SVCEN(20,20),
     2IDATA(6C0),NISV(20),VMEAN(20),SC(6C),XVEC(2C),YVEC(20),ET(20,20),
     3NISVT(20),SVCENT(20,20),IDATAT(6CC),VEC(20,20),EIG(20)
      DIMENSION SV(20,20),NI(20)
C
C     CALCULATE B AND THEN W = T - B
C     B = SUM  NIS(M)*SVCEN(M)**2
```

```
C
      DO 100 J=1,NVARS
  100 E(J,K) = 0.0
      DO 105 M=1,NGP
      DO 105 J=1,NVARS
      DO 105 K=J,NVARS
  105 BT(J,K) = SV(M,J) * SV(M,K)
      DO 110 J=1,NVARS
      DO 110 K=J,NVARS
      E(J,K) = E(J,K) + BT(J,K)* NI(M)
  110 W(J,K) = T(J,K) - E(J,K)
C
C     CALCULATE CHOLESKY FACTOR OF W IF ICIR = 2
C
      GO TO (120,115) ICIR
  115 DO 116 J=1,NVARS
      DO 116 K=J,NVARS
  116 WFCT(J,K) = W(J,K)
      CALL UPFFCT (NVARS,WFCT,M)
      IF (M.NE.0) GO TO 125
  120 RETURN
  125 IFINE=1
      WRITE (6,90C)
  9CC FORMAT (' THE WITHIN GROUPS MATRIX IS SINGULAR ')
      RETURN
      END
      FUNCTION URANC(IRANC)
C
C     THIS FUNCTION CALCULATES UNIFORMLY DISTRIBUTED RANDOM NUMBERS.
C     BETWEEN 0 AND 1
C  3**19 CONGRUENTIAL UNIFORM RANDOM NUMBER GENERATOR
C
      IRAND = IRAND*1162261467
    1 IF (IRAND.GT.0) GO TO 3
      IRAND = -IRAND
    3 URAND = FLOAT(IRAND)*0.4656612873E-9
      RETURN
      END
      SUBROUTINE LPRFCT(N,A,M)
C
C     REPLACE UPPER TRIANGLE OF A SQUARE POSITIVE DEFINITE MATRIX A
C     BY ITS CHOLESKI FACTOR
C
      DIMENSION A(20,20)
C     CLEAR THE ERROR INDICATOR
      M=0
      NJ=N-1
```

RRRC9540C
RRRC941C
RRRC9420C
RRRC9430C
RRRC944C
RRRC9450C
RRRC9460C
RRRC9470C
RRRC9480C
RRRC9490C
RRRC9500C
RRRC951C
RRRC9520C
RRRC9530C
RRRC9540C
RRRC955C
RRRC9560C
RRRC9570C
RRRC958C
RRRC9590C
RRRC9600C
RRRC9610C
RRRC962C
RRRC9630C
RRRC9640C
RRRC9650C
RRRC966C
RRRC9670C
RRRC9680C
RRRC9690C
RRRC9700C
RRRC9710C
RRRC9720C
RRRC973C
RRRC9740C
RRRC975C
RRRC9760C
RRRC9770C
RRRC9780C
RRRC9790C
RRRC9810C
RRRC982C
RRRC9830C
RRRC9840C
RRRC985C
RRRC9860C
RRRC987C

108

```
      IF(N1) 230,100,100
100   DC 220 K=1,N
      AKK=A(K,K)
      IF(K .EC. 1) GO TC 120
110   DC 110 J=2,K
      AKK=AKK-A(J-1,K)**2
120   IF(A(K,K)) 140,140,130
C     R IS MULT CCRRELATICN OF VARIABLE K WITH ALL FRECEEDING VARIABLES.
C     IN THE CASE OF A CCVARIANCE MATRIX AKK/A(K,K) IS 1-F**2 WHERE
130   IF(AKK/A(K,K) .GE. .001) GO TC 150
140   M=K
      AKK=0.
150   AKK=SCRT(AKK)
      A(K,K)=AKK
      IF(K .EC. N) GO TC 230
C
180   DC 220 I=K,N1
      AKI=A(K,I+1)
      IF(K .EC. 1) GO TC 190
190   DC 180 J=2,K
      AKI=AKI-A(J-1,K)*A(J-1,I+1)
200   IF(AKK) 210,200,210
210   A(K,I+1)=0.0
      GC TO 220
      A(K,I+1)=AKI/AKK
220   CCNTINUE
230   RETURN
      ENC
      SLBRCUTINE LTIRT(M,N,S,B)
C
C     INVERSE CF LPPER (S) TRANSPCSED TIMES RECTANGLE B TRANSFCSED.
C
      DIMENSICN S(20,20),B(1,20)
130   DCC 130 J=1,M
130   130 I=1,N
      SLM=0.0
50    IF (S(I,I)) 90,120,90
90    SLM=B(J,I)
110   IF(IM1) 120,120,100
100   CCC 110 K=1,IM1
110   SLM = SUM-S(K,I)*E(J,K)
130   B(J,I) = SUM/S(I,I)
      RETURN
      ENC
      SLBROUTINE LTISUI (N,A,B)
C
```

109

```
C     UPPER (B) TRANSPOSE INVERSE TIMES A TIMES UPPER (B) INVERSE.             RRRIC036C
C                                                                             RRRIC037C
      DIMENSION A(20,20),B(20,20)                                            RRRIC039C
C     NOTE THAT POSTMULT IS CARRIED OUT ON FINAL VALUES LEFT BY PREMULT      RRRIC040C
      DO 200 I=1,N                                                          RRRIC041C
      II=N-I                                                               RRRIC042C
C     NOTE THAT PREMULT ON RIGHT HALF OF ROW IS SAME AS POSTMULT            RRRIC043C
C     LOWER HALF OF COLUMN - EXCEPT FOR DIAG TERM WHICH IS THE              RRRIC044C
C     FINAL VALUE LEFT OF COLUMN - EXCEPT FOR DIAG TERM WHICH IS THE        RRRIC045C
      DO 130 J=1,N                                                         RRRIC046C
      DO 120 K=1,II                                                       RRRIC047C
 10C  IF(II) 11C,11C,11C                                                   RRRIC048C
 11CC A(J,I) = A(J,I) - B(K,I)*A(J,K)                                      RRRIC049C
 12CC A(J,I) = A(J,I)/B(I,I)                                              RRRIC050C
 13C  IF(B(I,I)) .EQ. 0.C) A(J,I) = 0.0                                    RRRIC051C
C     NOTE THAT ELEMENTS IN LEFT HALF OF ROW ARE FINAL FOR PREMULT         RRRIC052C
C     NOTE THAT DIAG ELEMENT WAS PREVIOUSLY THE FINAL RESULT OF            RRRIC053C
C     PREMULT. NOW WE MAKE THE FINAL RESULT OF POSTMULT.                    RRRIC054C
      DO 200 J=1,I                                                        RRRIC055C
      JI=I-J+1                                                            RRRIC056C
 14C  IF(JI) 16C,16C,14C                                                   RRRIC057C
C     HORIZONTAL BRANCH OF INNER PRODUCT EXCLUDING DIAG TERM               RRRIC058C
 15CC K=1,JI                                                               RRRIC059C
 15C  A(I,J)=A(I,J) - B(K,I)*A(J,K)                                        RRRIC060C
 16C  IF(II-J) 19C,170,170                                                 RRRIC061C
C     VERTICAL BRANCH OF INNER PRODUCT INCLUDING DIAG TERM                 RRRIC062C
 17CC K=J,II                                                               RRRIC063C
 18C  A(I,J)=A(I,J) - B(K,I)*A(K,J)                                        RRRIC064C
 19CC A(I,J)=A(I,J)/B(I,I)                                                 RRRIC065C
 20CC IF(B(I,I) .EQ. 0.0) A(I,J) = 0.0                                     RRRIC066C
      RETURN                                                              RRRIC067C
      END                                                                 RRRIC068C
      SUBROUTINE EIGN(NN,A,EIG,VEC,INC)                                   RRRIC069C
C     NN= SIZE OF MATRIX                                                  RRRIC070C
C     A = MATRIX (ONLY LOWER TRIANGLE IS USED + THIS IS DESTROYED)        RRRIC071C
C     EIG = RETURNED EIGENVALUES IN ALGEBRAIC DESCENDING ORDER            RRRIC072C
C     VEC = RETURNED EIGENVECTORS IN COLUMNS                              RRRIC073C
C     INC = ERROR RETURN INDICATOR                                        RRRIC074C
C         0 FOR NORMAL RETURN                                             RRRIC075C
C         1 SUM OF EIGENVALUES NOT EQUAL TO TRACE                         RRRIC076C
C         2 SUM OF EIGENVALUES SQUARED NOT EQUAL TO NORM                  RRRIC077C
C         3 BOTH OF THESE ERRORS                                          RRRIC078C
      DIMENSION A(20,20),VEC(20,20),GAMMA(20),EETA(20),BETASC(20),EIG(20) RRRIC079C
      DIMENSION W(20),Q(20)                                               RRRIC080C
C     THE FOLLOWING DIMENSIONED VARIABLES ARE EQUIVALENCED                RRRIC081C
      EQUIVALENCE (P(1),EETA(1)),(Q(1),EETA(1))                           RRRIC082C
      DIMENSION F(19),Q(19)                                               RRRIC083C
      DIMENSION IFOSV(20),IVPOS(20),ICRC(20)                              RRRIC083C
```

```
      EQUIVALENCE (IPOSV(1),GAMMA(1)),(IVPOS(1),BETA(1)),
     1(ICRD(1),BETASQ(1))
      N=NN
C     RESET ERROR RETURN INDICATOR
      IND=0
      IF(N .EQ. 0) GO TO 560
      N1=N-1
      N2=N-2
C     COMPUTE THE TRACE AND EUCLIDIAN NORM OF THE INPUT MATRIX
C     LATER CHECK AGAINST SUM AND SUM OF SQUARES OF EIGENVALUES
      TRACE=0.
      ENORM=0.
  100 DO 100 J=1,N
  100 DO 100 I=J,N
      ENORM=ENORM+A(I,J)**2
  110 TRACE=TRACE+A(J,J)
      ENORM=ENORM-.5*A(J,J)**2
      ENORM=ENORM+ENORM
      GAMMA(1)=A(1,1)
      IF(N2) 280,270,120
  120 DO 260 NR=1,N2
      EE=A(NR+1,NR)
      S=0.
  130 DO 130 I=NR,N2
      S=S+A(I+2,NR)**2
      A(NR+1,NR)=S
      IF(S) 250,250,140
  140 S=S+B*B
      SGN=+1
      IF(B) 150,160,160
  150 SGN=-1.
  160 SQRTS=SQRT(S)
      D=SGN/(SQRT(S+SQRTS)
      TEMP=SQRT(.5+B*C)
      W(NR)=TEMP
      A(NR+1,NR)=TEMP
      C=D/TEMP
C     B IS FACTOR OF PROPORTIONALITY. NOW COMPUTE AND SAVE W VECTOR.
C     EXTRA SINGLY SUBSCRIPTED W VECTOR USED FOR SPEED.
  170 DO 170 I=NR,N2
      TEMP=D*A(I+2,NR)
      W(I+1)=TEMP
      A(I+2,NR)=TEMP
C     PREMULTIPLY VECTOR W BY MATRIX A TO OBTAIN P VECTOR
C     SIMULTANEOUSLY ACCUMULATE DOT PRODUCT WP,(THE SCALAR K)
      WTAW=0.
```

                                                                    RRR10840
                                                                    RRR10850
                                                                    RRR10860
                                                                    RRR10870
                                                                    RRR10880
                                                                    RRR10890
                                                                    RRR10900
                                                                    RRR10910
                                                                    RRR10920
                                                                    RRR10930
                                                                    RRR10950
                                                                    RRR10960
                                                                    RRR10970
                                                                    RRR10980
                                                                    RRR10990
                                                                    RRR11000
                                                                    RRR11010
                                                                    RRR11020
                                                                    RRR11030
                                                                    RRR11040
                                                                    RRR11050
                                                                    RRR11060
                                                                    RRR11070
                                                                    RRR11080
                                                                    RRR11090
                                                                    RRR11100
                                                                    RRR11110
                                                                    RRR11120
                                                                    RRR11130
                                                                    RRR11140
                                                                    RRR11150
                                                                    RRR11160
                                                                    RRR11170
                                                                    RRR11180
                                                                    RRR11190
                                                                    RRR11200
                                                                    RRR11210
                                                                    RRR11220
                                                                    RRR11230
                                                                    RRR11240
                                                                    RRR11250
                                                                    RRR11260
                                                                    RRR11270
                                                                    RRR11280
                                                                    RRR11290
                                                                    RRR11300
                                                                    RRR11310

111

```
      DO 220 I=NR,N1
      SUM=0.  J=NR,I
180   SUM=SUM+A(I+1,J+1)*W(J)
      II=I+1
      IF(N1-I1) 210,150,190
150   DO 200 J=II,N1
200   SUM=SUM+A(J+1,I+1)*W(J)
210   P(I)=SUM
C     F VECTOR AND SCALAR K NOW STORED. NEXT COMPUTE Q VECTOR
220   SUM=SUM+TAU+SUM*W(I)
C
230   Q(I)=P(I)-W1AW*W(I)
C     NOW FORM PAF MATRIX, REQUIRED PART
      DO 240 J=NR,N1
      QJ=Q(J)
      WJ=W(J)
240   A(I+1,J+1)=A(I+1,J+1)-2.*(W(I)*QJ+WJ*Q(I))
250   BETA(NR)=B
      BETASQ(NR)=E*B
260   GAMMA(NF+1)=A(NR+1,NR+1)
270   EETA(N,N-1)=B
      EETA(N-1,1)=B
      GAMMA(N-1)=A(N,N)
      BETASQ(K-1)=B*B
280   BETASQ(N)=0.
C     ADJOIN AN IDENTITY MATRIX TO BE POSTMULTIPLIED BY ROTATIONS.
      DO 300 I=1,A
290   VEC(I,J)=0.
300   VEC(I,I)=1.
      M=N
      SUM=0.
      NPAS=1
      GO TO 400
310   SUM=SUM+SHIFT
      COSA=1.
      G=GAMMA(1)-SHIFT
      PP=G
      PPES=PP*PP+EETASQ(1)
      PPER=SQRT(PPES)
      DO 270 J=1,M
      COSAP=COSA
      IF(PPES .NE. 0.) GO TO 320
      SINA=0.
      SINA2=0.
      COSA=1.
```

```
320   GC TO 350                                                  RRR11800
      SINA=BETA(J)/PPBR                                           RRR11810
      SINA2=BETASQ(J)/PFBS                                        RRR11820
      CCSA=FP/PPBR                                                RRR11830
C     PCSTMULTIPLY IDENTITY BY P-TRANSPCSE MATRIX                 RRR11840
      NT=J+NPAS                                                   RRR11850
      IF(NT .GE. N) NT=N                                          RRR11860
330   CC 340 I=1,NT                                               RRR11870
      TEMP=COSA*VEC(I,J)+SINA*VEC(I,J+1)                          RRR11880
      VEC(I,J+1)=-SINA*VEC(I,J)+CCSA*VEC(I,J+1)                   RRR11890
340   VEC(I,J)=TEMP                                               RRR11900
350   CIA=GAMMA(J+1)-SHIFT                                        RRR11910
      L=SINA2*(G+CIA)                                             RRR11920
      GAMMA(J)=G+L                                                RRR11930
      G=CIA-U                                                     RRR11940
      PF=DIA*COSA-SINA*CCSAP*BETA(J)                              RRR11950
      IF(J .NE. M) GO TC 360                                      RRR11960
      BETA(J)=SINA*PP                                             RRR11970
      BETASQ(J)=SINA2*PF*PP                                       RRR11980
      GC TO 380                                                   RRR11990
360   FFES=PP*PP+EETASQ(J+1)                                      RRR12000
      PFBR=SQRT(PFBS)                                             RRR12010
      EETA(J)=SINA*PPBR                                           RRR12020
      BETASQ(J)=SINA2*PPBS                                        RRR12030
380   GAMMA(M+1)=C                                                RRR12040
C     TEST FOR CONVERGEACE CF LAST DIAGCNAL ELEMENT               RRR12050
      NFAS=NPAS+1                                                 RRR12060
      IF(BETASQ(M).GT..1.E-21) GC TO 41C                          RRF12070
      EIG(M+1)=GAMMA(M+1)+SUM                                     RRR12080
350   BETA(M)=0.                                                  RRR12090
4CC   BETASQ(M)=C.                                                RRR12100
      M=M-1                                                       RRR12110
      IF(M .EQ. C) GO TC 43C                                      RRR12120
      IF(EETASQ(M) .LE. 1.E-21) GO TO 350                         RRR12130
C     TAKE ROCT OF CORNER 2 BY 2 NEAREST TO LOWER CIAGONAL IN VALUE RRR12140
C     AS ESTIMATE OF EIGENVALUE TO USE FCR SHIFT                  RRR12150
41C   A2=GAMMA(M+1)                                               RRR12160
      R2=.5*A2                                                    RRR12170
      R12=.5*GAMMA(M)                                             RRR12180
      CIF2=R1+R2                                                  RRR12190
      CIF=R1-R2                                                   RRR12200
      TEMP=SQRT(DIF*DIF+BETASQ(M))                                RRR12210
      R1=R12+TEMP                                                 RRR12220
      R2=R12-TEMP                                                 RRR12230
      CIF=ABS(A2-R1)-ABS(A2-R2)                                   RRR12240
      IF(DIF.LT. 0.) GC TO 420                                    RRR12250
      SHIFT=R2                                                    RRR12260
      GC TO 310                                                   RRR12270
```

113

```
420   SHIFT=R1                                                      RRR1228C
430   GO TO 310                                                     RRR1229C
      EIG(I)=GAMMA(I)+SLM                                           RRR1230C
C     INITIALIZE AUXILIARY TABLES REQUIRED FOR REARRANGING THE VECTORS  RRR1231C
      DO 440 J=1,N                                                  RRR1232C
      IPCSV(J)=J                                                    RRR1233C
      IVPOS(J)=J                                                    RRR1234C
440   ICRC(J)=J                                                     RRR1235C
C     USE A TRANSPOSITION SORT TO ORDER THE EIGENVALUES             RRR1236C
      M=N                                                           RRR1237C
450   GO TO 470                                                     RRR1238C
460   DO 460 J=1,M                                                  RRR1239C
      IF(EIG(J).GE. EIG(J+1)) GO TO 460                             RRR1240C
      TEMP=EIG(J)                                                   RRR1241C
      EIG(J)=EIG(J+1)                                               RRR1242C
      EIG(J+1)=TEMP                                                 RRR1243C
      ITEMP=ICRC(J)                                                 RRR1244C
      ICRC(J)=ICRC(J+1)                                             RRR1245C
      ICRC(J+1)=ITEMP                                               RRR1246C
460   CONTINUE                                                      RRR1247C
470   M=M-1                                                         RRR1248C
      IF(M .NE. 0) GO TO 450                                        RRR1249C
      IF(N1 .EQ. 0) GO TO 500                                       RRR1250C
490   DO 490 L=1,N1                                                 RRR1251C
      NV=IORO(L)                                                    RRR1252C
      NF=IPCSV(NV)                                                  RRR1253C
      IF(NP .EQ. L) GO TO 490                                       RRR1254C
      LV=IVPOS(L)                                                   RRR1255C
      INPCS(NP)=L                                                   RRR1256C
      IPOSV(LV)=NF                                                  RRR1257C
480   DO 480 I=1,N                                                  RRR1258C
      TEMP=VEC(I,L)                                                 RRR1259C
      VEC(I,L)=VEC(I,NF)                                            RRR1260C
      VEC(I,NF)=TEMP                                                RRR1261C
490   CONTINUE                                                      RRR1262C
500   SLM=0.                                                        RRR1263C
C     BACK TRANSFORM THE VECTORS OF THE TRIPLE DIAGONAL MATRIX      RRR1264C
      DO 550 NRR=1,N                                                RRR1265C
510   K=N1                                                          RRR1266C
      K=K-1                                                         RRR1267C
      IF(K .LE. 0) GO TO 540                                        RRR1268C
      SUM=0.                                                        RRR1269C
520   DO 520 I=K,N1                                                 RRR1270C
      SUM=SUM+VEC(I+1,NRR)*A(I+1,K)                                 RRR1271C
      SUM=SLM+SUM                                                   RRR1272C
530   DO 530 I=K,N1                                                 RRR1273C
      VEC(I+1,NRR)=VEC(I+1,NRR)-SUM*A(I+1,K)                        RRR1274C
                                                                    RRR1275C
```

114

```
540   CC TO 51C
      ESCM=ESCM+EIG(NRR)                                      RRR1276C
55C   ESSQ=ESSQ+EIG(NRR)**2                                   RRR1277C
      TEMP=ABS(51C*TRACE)                                     RRR12780
      IF((ABS(TRACE-ESUM)+TEMP)-TEMP .NE. 0.) INC=INC+1       RRR12790
      TEMP=1024.*ENORM                                        RRR12800
      IF((ABS(ENORM-ESSC)+TEMP)-TEMP .NE. C.) INC=INC+2       RRF12810
56C   RETLRN                                                  RRF12820
      EAC                                                     RRR1283C
                                                              RRR12840
```

115

Appendix  C:

```
RUN NAME        PRINCIPAL CCMPCNENT ANALYSIS OF TANKS
VARIABLE LIST   NATIONS,X1 TC X10
INPUT MEDIUM    CARD
INPUT FORMAT    FIXED(1X,F2.0,5X,8F8.3/8X,2F8.3)
N OF CASES      24
VAR LABELS      NATIONS TYPE OF TANK/
                X1 WEIGHT/X2 LENGTH/X3 WIDTH/X4 HEIGHT/X5 RCAD SPEED/
                X6 TRENCH CRCSSING/X7 GRCUNC PRESSLRE/X8 MAX ARMAMENT/
                X9 GROUND CLEARANCE/X10 FCWER TO ENGINE RATIC/
                VARIABLES=X1 TO X10/TYPE=PA1/NFACTCRS=10/
                8,11
FACTOR          ALL
OPTICNS
STATISTICS
READ INPUT DATA
          CECK
          . . . . .

FINISH
```

116

```
RUN NAME        TANK DISCRIMINANT
VARIABLE LIST   NATIONS,X1 TO X10
NO OF CASES     24
INPUT MEDIUM    CARD
INPUT FORMAT    FIXED(1X,F2.0,5X,8F8.3/6X,2F8.3)
RECODE          NATIONS(11 THRU 13=1)(15,16,17,24,25,26,30,31,32,33=2)
                (19,20,21,22,23,28,29,=3)(14,18,27,34=4)
VAR LABELS      NATIONS TYPE OF TANK/
                X1 WEIGHT/X2 LENGTH/X3 WIDTH/X4 HEIGHT/X5 ROAD SPEED/
                X6 TRENCH CROSSING/X7 GROUND PRESSURE/X8 MAX ARMAMENT/
                X9 GROUND CLEARANCE/X10 POWER TO ENGINE RATIO/
DISCRIMINANT    GROUPS=NATIONS(1,4)/VARIABLES=X1 TO X10/
                ANALYSIS=X1 TO X10/METHOD=MAHAL/
OPTIONS         5,6,7,9,11,12
STATISTICS      ALL
READ INPUT DATA
      DECK
      .  .  .  .
FINISH
```

# BIBLIOGRAPHY

1.  Aiken, J. W., "Development of Cluster Analysis for Student Opinion Data", Master's Thesis, Naval Post-graduate School, Monterey, CA, 1979.

2.  Anderson, T. W., An Introduction to Multivariate Statistical Analysis, Wiley, 1958.

3.  Anderberg, M. R., Cluster Analysis for Applications, Academic Press, 1973.

4.  Barr, D. R., and Richards, F. R., Utility Assessment Methodology (Report on Relative Utility Score of the AN-TPQ/27). System Exploration Inc., Monterey, CA, 1980.

5.  Eisenheis, R. A., and Avery, R. B., Discriminant Analysis and Classification Procedures, Lexington Books, 1972.

6.  Friedman, H. P., and Rubin, J., "On Some Invariant Criteria for Grouping Data", Journal of the American Statistical Association, Vol. 62: 320, Dec. 1967.

7.  Giri, N. C., Multivariate Statistical Inference, Academic Press, 1977.

8.  Gnanadesikan, R., Methods for Statistical Data Analysis of Multivariate Observations, Wiley, 1977.

9.  Green, P. E., and Carroll, J. D., Mathematical Tools for Applied Multivariate Analysis, Academic Press, 1976.

10. Hartigan, J. A., Clustering Algorithms, Wiley, 1975.

11. Kandell, M. G., The Advanced Theory of Statistics, Vol. III, Charles Griffin and Company, 1947.

12. Keeney, R. L., and Raiffa, H., Decisions with Multiple Objectives, Wiley, 1976.

13. Klecka, W. R., "Discriminant Analysis", SPSS (Statistical Package for the Social Sciences), McGraw Hill, 1975.

14. Kim, J. O., "Factor Analysis", SPSS (Statistical Package for the Social Sciences), McGral Hill, 1975.

15. MacQueen, J., "Some Methods for Classification and Analysis of Multivariate Observations", Paper presented at Fifth Berkeley Symposium on Mathematical Statistics and Probability, University of California, 1965-1966.

16. McRae, D. J., Clustering Multivariate Observations, Doctoral Dissertation, University of North Carolina, Chapel Hill, N.C., 1973.

17. Morrison, D. F., Multivariate Analysis: Technique for Educational and Psychological Research, Wiley, 1971.

18. The IMSL Library, Vol. 3, IMSL, Inc., 1979.

INITIAL DISTRIBUTION LIST

|  |  | No. Copies |
|---|---|---|
| 1. | Defense Technical Information Center (DTIC)<br>Cameron Station<br>Alexandria, Virginia 22314 | 2 |
| 2. | Library, Code 0142<br>Naval Postgraduate School<br>Monterey, California 93940 | 2 |
| 3. | Department Chairman, Code 55<br>Department of Operations Research<br>Naval Postgraduate School<br>Monterey, California 93940 | 4 |
| 4. | Professor F. R. Richards, Code 55Rh<br>(Thesis Advisor)<br>Department of Operations Research<br>Naval Postgraduate School<br>Monterey, California 93940 | 2 |
| 5. | Professor D. R. Barr, Code 55Bn<br>(Second Reader)<br>Department of Operations Research<br>Naval Postgraduate School<br>Monterey, California 93940 | 1 |
| 6. | LTC Jin Ki Lee, ROK Army (Student)<br>6-60$^1$, Gun-In Apt, Dong-Bing-Go Dong<br>Yong-San, Seoul<br>Korea | 5 |